

Supplementary Information

Supplementary Note 1: SBayesS model

Let us consider an individual-level data-based multiple regression model in a GWAS data set:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e} \quad (1)$$

where \mathbf{y} is the vector of phenotypes adjusted for all fixed effects, \mathbf{X} is the column-centered genotype matrix, $\boldsymbol{\beta}$ is the vector of SNP effects, and \mathbf{e} is the vector of residuals with $Var(\mathbf{e}) = \mathbf{I}\sigma_e^2$ for a sample of unrelated individuals. Assuming Hardy-Weinberg equilibrium (HWE), the variance of the genotypes of SNP j is $h_j = 2p_jq_j$, where p_j is the minor allele frequency (MAF) and $q_j = 1 - p_j$. Let \mathbf{D} be a diagonal matrix with $D_j = \mathbf{X}'_j\mathbf{X}_j = h_jn_j$, where n_j is per-SNP sample size. Multiplying both sides of (1) by $\mathbf{D}^{-1}\mathbf{X}'$ gives

$$\mathbf{D}^{-1}\mathbf{X}'\mathbf{y} = \mathbf{D}^{-1}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} + \mathbf{D}^{-1}\mathbf{X}'\mathbf{e} \quad (2)$$

Note that $\mathbf{D}^{-1}\mathbf{X}'\mathbf{y} = \mathbf{b}$, the vector of least square estimates of SNP marginal effects from GWAS, and $\mathbf{D}^{-1}\mathbf{X}'\mathbf{X} = \mathbf{D}^{-\frac{1}{2}}\mathbf{B}\mathbf{D}^{\frac{1}{2}}$ where $\mathbf{B} = \mathbf{D}^{-\frac{1}{2}}\mathbf{X}'\mathbf{X}\mathbf{D}^{-\frac{1}{2}}$ is the linkage disequilibrium (LD) correlation matrix among all SNPs (3). Let $\boldsymbol{\epsilon} = \mathbf{D}^{-1}\mathbf{X}'\mathbf{e}$. Then, (2) can be written as

$$\mathbf{b} = \mathbf{D}^{-\frac{1}{2}}\mathbf{B}\mathbf{D}^{\frac{1}{2}}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad (3)$$

Or, in a scalar form,

$$b_j = \sum_{k=1}^m \sqrt{\frac{h_k n_k}{h_j n_j}} B_{jk} \beta_k + \epsilon_j \quad (4)$$

with m being the total number of SNPs. Let $\sigma_{X_j}^2$, $\sigma_{X_k}^2$ and σ_{X_j, X_k} denote the genotype variance of SNP j and k and their covariance. Then, (4) can be simplified to

$$b_j = \sum_{k=1}^m \sqrt{\frac{n_k}{n_j}} \beta_{X_j, X_k} \beta_k + \epsilon_j \quad (5)$$

where $\beta_{X_j, X_k} = \sigma_{X_j, X_k} / \sigma_{X_j}^2$ is the regression of SNP k on that of SNP j . In other words, we model the marginal effect of each SNP as a linear combination of other SNP effects with the weights being a function of regression coefficient of SNP genotypes and per-SNP sample sizes. In contrast to the identity structure of residual variance in (1), the residuals in (3) are not independent in the presence of LD, because

$$\begin{aligned} Var(\boldsymbol{\epsilon}) &= Var(\mathbf{D}^{-1}\mathbf{X}'\mathbf{e}) \\ &= \mathbf{D}^{-1}\mathbf{X}'\mathbf{X}\mathbf{D}^{-1}\sigma_e^2 \\ &= \mathbf{D}^{-\frac{1}{2}}\mathbf{B}\mathbf{D}^{-\frac{1}{2}}\sigma_e^2 \end{aligned} \quad (6)$$

Let $\mathbf{W} = \mathbf{D}^{-\frac{1}{2}}\mathbf{B}\mathbf{D}^{\frac{1}{2}}$ and $\mathbf{R} = \mathbf{D}^{-\frac{1}{2}}\mathbf{B}\mathbf{D}^{-\frac{1}{2}}$. Finally, from (3) we have

$$\mathbf{b} = \mathbf{W}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad (7)$$

with $Var(\boldsymbol{\epsilon}) = \mathbf{R}\sigma_e^2$. This is a generic form of summary-data-based Bayesian regressions (SBR), which is similar to the RSS model of Zhu and Stephens [2].

As in BayesS [3], we assume the effect size is related to MAF through parameter S :

$$\beta_j \begin{cases} \sim N(0, h_j^S \sigma_\beta^2), & \pi \\ = 0, & 1 - \pi \end{cases} \quad (8)$$

where S , σ_β^2 and π are considered as unknown. The prior for S is a standard normal distribution

$$S \sim N(0, 1)$$

The prior for π is a uniform between zero and one, namely

$$\pi \sim Beta(1, 1)$$

The prior for σ_β^2 is a scaled inverse chi-square distribution

$$\sigma_\beta^2 \sim \nu_\beta \tau_\beta^2 \chi_{\nu_\beta}^{-2}$$

where $\nu_\beta = 4$ and

$$\tau_\beta^2 = \frac{\nu_\beta - 2}{\nu_\beta} \frac{V_P h_0^2}{\pi_0 \sum_j h_j^{S_0+1}}$$

where V_P is the phenotypic variance estimated from the summary statistics (as shown below) and h_0^2 , π_0 and S_0 are the prior knowledge of SNP-based heritability, π and S , respectively. Similarly, we give a scaled inverse chi-square prior for σ_e^2 in (6)

$$\sigma_e^2 \sim \nu_e \tau_e^2 \chi_{\nu_e}^{-2}$$

where $\nu_e = 4$ and

$$\tau_e^2 = \frac{\nu_e - 2}{\nu_e} V_P (1 - h_0^2)$$

We estimate the phenotypic variance of the trait following Yang et al [1]'s approach, which is based on the standard error of the marginal SNP effect estimate from GWAS. Because

$$SE_j^2 = \frac{\sigma_j^2}{\mathbf{X}_j' \mathbf{X}_j} = \frac{\mathbf{y}' \mathbf{y} - \mathbf{X}_j' \mathbf{X}_j b_j^2}{\mathbf{X}_j' \mathbf{X}_j}$$

where σ_j^2 is the residual variance for the GWAS model fitting SNP j . Rearranging gives

$$V_{P,j} = \frac{\mathbf{y}'\mathbf{y}}{n_j} = D_{jj} \left(SE_j^2 + \frac{b_j^2}{n_j} \right) \quad (9)$$

The phenotypic variance V_P is then calculated as the median of $V_{P,j}$ across all SNPs [1]. Since $D_{jj} = 2p_jq_jn_j$ but the allele frequencies from the publicly available summary data are often not exact, we substitute $V_{P,j}$ by V_P in (9) to reestimate p_j given the input values of SE_j , b_j and n_j .

We call model (7) with the above prior as ‘‘SBayesS’’. Specifying a different prior distribution to β_j gives SBR form of other Bayesian alphabet models. For example, a mixture prior of normals with different variances for β_j under SBR framework becomes SBayesR [4].

As shown above, when the LD correlations are computed using all SNPs in the GWAS sample, SBayesS model is a linear transformation of the BayesS model without loss of information, in which case the two models are equivalent in terms of posterior inference (Section). However, it is impractical to store pairwise LD correlations of all genome-wide SNPs in the computer memory and not always feasible to access individual-level genotypes of the GWAS sample. Thus, we propose to use a sparse LD matrix that is computed from a reference sample, ignoring the small LD correlation estimates due to sampling variation. In this case, SBayesS becomes an approximation to BayesS. Assuming the LD reference sample is a random draw from the same population of the GWAS sample, the discrepancy between SBayesS and BayesS arises from the sampling variance of LD correlations used in SBayesS. Ignoring the sampling variance of LD estimates may cause a failure to converge in the Markov chain Monte Carlo (MCMC) sampling process or a bias in parameter estimation (Section). In this study, we model analytically the sampling variance of LD estimates as part of the residual variance and allow the estimate of residual variance to vary across SNPs (Section). We show the MCMC sampling scheme for the model parameters in Section and an efficient updating strategy in MCMC in Section .

Equivalence between SBayesS and BayesS

Here, we show that when the LD correlations are computed using all SNPs in the GWAS sample, SBayesS and BayesS models are equivalent in terms of posterior inference. Without loss of generality we assume $\pi = 0$, the posterior distribution of β in SBayesS is

$$\begin{aligned} f(\beta | \mathbf{b}, else) &\propto f(\mathbf{b} | \beta, else) f(\beta) \\ &\propto \exp \left\{ -\frac{1}{2\sigma_e^2} (\mathbf{b} - \mathbf{W}\beta)' \mathbf{R}^{-1} (\mathbf{b} - \mathbf{W}\beta) \right\} \exp \left\{ -\frac{\beta' \mathbf{G}^{-1} \beta}{2\sigma_\beta^2} \right\} \\ &= \exp \left\{ -\frac{1}{2\sigma_e^2} \left[\mathbf{b}' \mathbf{R}^{-1} \mathbf{b} - 2\beta' \mathbf{W}' \mathbf{R}^{-1} \mathbf{b} + \beta' \left(\mathbf{W}' \mathbf{R}^{-1} \mathbf{W} + \frac{\sigma_e^2}{\sigma_\beta^2} \mathbf{G}^{-1} \right) \beta \right] \right\} \\ &\propto \exp \left\{ -\frac{1}{2\sigma_e^2} \left[\beta' \left(\mathbf{W}' \mathbf{R}^{-1} \mathbf{W} + \frac{\sigma_e^2}{\sigma_\beta^2} \mathbf{G}^{-1} \right) \beta + 2\beta' \mathbf{W}' \mathbf{R}^{-1} \mathbf{b} \right] \right\} \end{aligned}$$

Note that

$$\begin{aligned}
\mathbf{W}'\mathbf{R}^{-1}\mathbf{W} &= \mathbf{D}^{\frac{1}{2}}\mathbf{B}\mathbf{D}^{-\frac{1}{2}}\mathbf{D}^{\frac{1}{2}}\mathbf{B}^{-1}\mathbf{D}^{\frac{1}{2}}\mathbf{D}^{-\frac{1}{2}}\mathbf{B}\mathbf{D}^{\frac{1}{2}} \\
&= \mathbf{D}^{\frac{1}{2}}\mathbf{B}\mathbf{D}^{\frac{1}{2}} \\
&= \mathbf{D}^{\frac{1}{2}}\mathbf{D}^{-\frac{1}{2}}\mathbf{X}'\mathbf{X}\mathbf{D}^{-\frac{1}{2}}\mathbf{D}^{\frac{1}{2}} \\
&= \mathbf{X}'\mathbf{X}
\end{aligned}$$

and

$$\begin{aligned}
\mathbf{W}'\mathbf{R}^{-1}\mathbf{b} &= \mathbf{D}^{\frac{1}{2}}\mathbf{B}\mathbf{D}^{-\frac{1}{2}}\mathbf{D}^{\frac{1}{2}}\mathbf{B}^{-1}\mathbf{D}^{\frac{1}{2}}\mathbf{b} \\
&= \mathbf{D}\mathbf{b} \\
&= \mathbf{X}'\mathbf{y}
\end{aligned}$$

Thus, the posterior distribution becomes

$$f(\boldsymbol{\beta}|\mathbf{b}, else) \propto \exp\left\{-\frac{1}{2\sigma_e^2}\left[\boldsymbol{\beta}'\left(\mathbf{X}'\mathbf{X} + \frac{\sigma_e^2}{\sigma_\beta^2}\mathbf{G}^{-1}\right)\boldsymbol{\beta} + 2\boldsymbol{\beta}'\mathbf{X}'\mathbf{y}\right]\right\}$$

This is equivalent to the posterior distribution given the individual genotype and phenotype data. It can be shown that the above is the kernel of a multivariate normal distribution, i.e.

$$\boldsymbol{\beta}|\mathbf{b}, else \sim MVN(\mathbf{C}^{-1}\mathbf{r}, \mathbf{C}^{-1}\sigma_e^2)$$

where $\mathbf{C} = \mathbf{W}'\mathbf{R}^{-1}\mathbf{W} + \frac{\sigma_e^2}{\sigma_\beta^2}\mathbf{G}^{-1} = \mathbf{X}'\mathbf{X} + \frac{\sigma_e^2}{\sigma_\beta^2}\mathbf{G}^{-1}$ and $\mathbf{r} = \mathbf{W}'\mathbf{R}^{-1}\mathbf{b} = \mathbf{X}'\mathbf{y}$. It is recognized that \mathbf{C} and \mathbf{r} are the left- and right-hand sides of the mixed-model equations

$$\underbrace{[\mathbf{W}'\mathbf{R}^{-1}\mathbf{W} + \mathbf{G}^{-1}\lambda]}_{\mathbf{C}}\boldsymbol{\beta} = \underbrace{\mathbf{W}'\mathbf{R}^{-1}\mathbf{b}}_{\mathbf{r}} \quad (10)$$

with $\lambda = \sigma_e^2/\sigma_\beta^2$. It can be further shown that in the Gibbs sampling, the full conditional distribution of β_j is

$$\beta_j|\mathbf{b}, \boldsymbol{\beta}_{-j}, else \sim N\left(\frac{r_j}{C_j}, \frac{\sigma_e^2}{C_j}\right) \quad (11)$$

where

$$\begin{aligned}
r_j &= D_j b_j - \sum_{k \neq j} D_j^{\frac{1}{2}} B_{jk} D_k^{\frac{1}{2}} \beta_k \\
&= \mathbf{X}'_j \mathbf{y} - \sum_{k \neq j} \mathbf{X}'_j \mathbf{X}_k \beta_k \\
C_j &= D_j + \frac{\sigma_e^2}{h_j^S \sigma_\beta^2} = \mathbf{X}'_j \mathbf{X}_j + \frac{\sigma_e^2}{h_j^S \sigma_\beta^2}
\end{aligned} \quad (12)$$

Consequence of ignoring LD sampling variation

The consequence of using a sparse LD matrix is that it can potentially bias the mean of the full conditional distribution for β_j . Let k denote a SNP in nonzero LD with the target SNP and l denote a SNP that does not have significant LD with the target SNP. The r_j in (12) can be written as

$$\begin{aligned}
 r_j &= D_j b_j - \sum_k D_j^{\frac{1}{2}} B_{jk} D_k^{\frac{1}{2}} \beta_k - \sum_l D_j^{\frac{1}{2}} B_{jl} D_l^{\frac{1}{2}} \beta_l \\
 &= D_j b_j - \sum_k D_j^{\frac{1}{2}} B_{jk} D_k^{\frac{1}{2}} \beta_k - t_j \\
 &= r_j^* - t_j
 \end{aligned} \tag{13}$$

The full conditional of β_j in (11) becomes

$$\beta_j | t_j \sim N \left(\frac{r_j^* - t_j}{C_j}, \frac{\sigma_e^2}{C_j} \right)$$

When reduced LD matrix is used, B_{jl} is set to be zero, we therefore completely ignore t_j and do not adjust the mean of the full conditional by the effects of SNPs in very low LD. However, although individual LD is trivial, the sum of them can be substantial after multiplied by n because $Var(\sum_l \mathbf{X}_j' \mathbf{X}_l) = D_j Var(\sum_l B_{jl}) \propto nm$ under the null. This may break the property of MCMC and fail the Gibbs sampling.

Modelling LD sampling variance

The use of a sparse LD matrix from a reference sample will result in two sources of sampling variation. The first is the difference in sampling variance between the reference and GWAS samples for the LD correlations included in the sparse LD matrix. The second is the sampling variance of LD correlations that are set to be zero. As shown above, ignoring these sampling variations will result in a bias of the mean in the full conditional distribution of β_j and thereby biases in the estimation of other model parameters. Here, we account for both sampling variations in the model, as described below.

Suppose the observed LD correlation between SNP j and k equal to the true population LD (ρ_{jk}) plus a deviation (tilde refers to the LD reference sample):

$$\begin{aligned}
 B_{jk} &= \rho_{jk} + \delta_{jk} \\
 \tilde{B}_{jk} &= \rho_{jk} + \tilde{\delta}_{jk}
 \end{aligned}$$

Then, the LD correlation in the GWAS sample is

$$B_{jk} = \begin{cases} \tilde{B}_{jk} + (\delta_{jk} - \tilde{\delta}_{jk}) & \text{if } \rho_{jk} \neq 0 \\ \delta_{jk} & \text{if } \rho_{jk} = 0 \end{cases} \tag{14}$$

Let Δ_{jk} denote the unobserved quantity in (14), i.e.

$$\Delta_{jk} = \begin{cases} \delta_{jk} - \tilde{\delta}_{jk} & \text{if } \rho_{jk} \neq 0 \\ \delta_{jk} & \text{if } \rho_{jk} = 0 \end{cases}$$

In (7), we can write

$$\mathbf{W} = \mathbf{W}^+ + \mathbf{W}^- \quad (15)$$

where $\mathbf{W}^+ = \mathbf{D}^{-\frac{1}{2}} \tilde{\mathbf{B}} \mathbf{D}^{\frac{1}{2}}$ is the observed data, and $\mathbf{W}^- = \mathbf{D}^{-\frac{1}{2}} \mathbf{\Delta} \mathbf{D}^{\frac{1}{2}}$ is not observed. Substituting (15) in (7) give

$$\mathbf{b} = \mathbf{W}^+ \boldsymbol{\beta} + \boldsymbol{\eta} \quad (16)$$

where $\boldsymbol{\eta} = \mathbf{W}^- \boldsymbol{\beta} + \boldsymbol{\epsilon}$ are the new residuals that contain difference in sampling deviation of LD between the GWAS and reference samples when the population LD is not zero, and sampling deviation of LD in the GWAS sample when the population LD is zero.

Conditional on $\mathbf{\Delta}$, the residual variance is

$$\text{Var}(\boldsymbol{\eta} | \mathbf{\Delta}) = \mathbf{W}^{-\prime} \mathbf{G} \mathbf{W}^- \sigma_{\beta}^2 + \mathbf{R} \sigma_e^2 \quad (17)$$

with $\mathbf{R} = \mathbf{D}^{-\frac{1}{2}} \tilde{\mathbf{B}} \mathbf{D}^{-\frac{1}{2}}$. However, this cannot be computed because $\mathbf{\Delta}$ is not observed.

Considering both $\boldsymbol{\beta}$ and $\mathbf{\Delta}$ as random, the diagonal value of $\text{Var}(\boldsymbol{\eta})$ in 16 is

$$\begin{aligned} \text{Var}(\eta_j) &= E[\text{Var}(\eta_j | \mathbf{\Delta}_j)] + \text{Var}[E(\eta_j | \mathbf{\Delta}_j)] \\ &= E\left[(\mathbf{W}_j^-)^\prime \mathbf{G} \mathbf{W}_j^- \sigma_{\beta}^2 \pi + D_j^{-1} \sigma_e^2\right] + 0 \\ &= E\left[\sum_{k=1}^m D_j^{-1} \Delta_{jk}^2 D_k G_k \sigma_{\beta}^2 \pi\right] + D_j^{-1} \sigma_e^2 \\ &= D_j^{-1} \sigma_{\beta}^2 \pi E[D_k G_k] E\left[\sum_k \Delta_{jk}^2\right] + D_j^{-1} \sigma_e^2 \\ &= D_j^{-1} \sigma_{\beta}^2 \pi \frac{\sum_k (2p_k q_k)^{S+1} n_j}{m} E\left[\sum_k \Delta_{jk}^2\right] + D_j^{-1} \sigma_e^2 \\ &= D_j^{-1} \left(\pi \sum_k (2p_k q_k)^{S+1} \sigma_{\beta}^2\right) \frac{n_j}{m} E\left[\sum_k \Delta_{jk}^2\right] + D_j^{-1} \sigma_e^2 \\ &= D_j^{-1} \sigma_g^2 \frac{n_j}{m} E\left[\sum_k \Delta_{jk}^2\right] + D_j^{-1} \sigma_e^2 \\ &= D_j^{-1} \left(\frac{n_j}{m} \sum_k E[\Delta_{jk}^2] \sigma_g^2 + \sigma_e^2\right) \\ &= D_j^{-1} \left(\frac{n_j}{m} \sum_k \text{Var}(\Delta_{jk}) \sigma_g^2 + \sigma_e^2\right) \end{aligned} \quad (18)$$

When $\rho_{jk} \neq 0$,

$$\begin{aligned} \text{Var}(\Delta_{jk}) &= \text{Var}(\delta_{jk}) + \text{Var}(\tilde{\delta}_{jk}) - 2\text{Cov}(\delta_{jk}, \tilde{\delta}_{jk}) \\ &= \frac{(1 - \rho_{jk}^2)^2}{n_j} + \frac{(1 - \rho_{jk}^2)^2}{\tilde{n}_j} - 2\text{Cov}(\delta_{jk}, \tilde{\delta}_{jk}) \end{aligned} \quad (19)$$

If there is no sample overlap between the LD reference and GWAS samples, then $\text{Cov}(\delta_{jk}, \tilde{\delta}_{jk}) = 0$. If we compute LD from the GWAS sample itself, then $\text{Var}(\Delta_{jk}) = 0$.

When $\rho_{jk} = 0$,

$$\text{Var}(\Delta_{jk}) = \text{Var}(\delta_{jk}) = \frac{1}{n_j}$$

If we estimate ρ_{jk} by \tilde{B}_{jk} and assume there is no sample overlap between the LD reference and GWAS samples, substituting these results in (18) gives

$$\text{Var}(\eta_j) = D_j^{-1} \left[\left(\frac{n_j}{m} s_j^2 + \frac{m_j^0}{m} \right) \sigma_g^2 + \sigma_e^2 \right] \quad (20)$$

where m_j^0 is the number of SNPs not in population LD with SNP j , σ_g^2 is the trait genetic variance, and

$$s_j^2 = \sum_{k=1}^{m_j} \left[\frac{(1 - \rho_{jk}^2)^2}{n_j} + \frac{(1 - \rho_{jk}^2)^2}{\tilde{n}_j} - 2\text{Cov}(\delta_{jk}, \tilde{\delta}_{jk}) \right]$$

is the total sampling variance for non-zero LD. In practice, we approximate ρ_{jk} by \tilde{B}_{jk} . In the absence of sample overlap between the LD reference and GWAS samples, $\text{Cov}(\delta_{jk}, \tilde{\delta}_{jk}) = 0$. In the case of complete sample overlap, $s_j^2 = 0$.

Similarly, considering both β and Δ as random, the off-diagonal value of $\text{Var}(\eta)$ in 16 is

$$\begin{aligned} \text{Cov}(\eta_j, \eta_k) &= E[\text{Cov}(\eta_j, \eta_k | \Delta_j, \Delta_k)] + \text{Cov}[E(\eta_j | \Delta_j), E(\eta_k | \Delta_k)] \\ &= E \left[(\mathbf{W}_j^-)' \mathbf{G} \mathbf{W}_k^- \sigma_\beta^2 + D_j^{-\frac{1}{2}} \tilde{B}_{jk} D_k^{-\frac{1}{2}} \sigma_e^2 \right] + 0 \\ &= E \left[\sum_l D_j^{-\frac{1}{2}} D_k^{-\frac{1}{2}} \Delta_{jl} \Delta_{kl} D_l G_l \sigma_\beta^2 \pi \right] + D_j^{-\frac{1}{2}} \tilde{B}_{jk} D_k^{-\frac{1}{2}} \sigma_e^2 \\ &= 0 + D_j^{-\frac{1}{2}} \tilde{B}_{jk} D_k^{-\frac{1}{2}} \sigma_e^2 \\ &= R_{jk} \sigma_e^2 \end{aligned}$$

Thus, the sampling variance of LD correlations only affect the diagonal but not off-diagonal values of the residual variance.

According to the derivation above, we have the following observations:

1. The LD sampling variance only affects the variance but not covariance of the model residuals. Thus, accounting for the LD sampling variance in the Gibbs sampling of β_j is straightforward (see below).

2. The LD sampling variation has two components, one due to the use of a different reference sample for LD information and the other due to the use of a sparse LD matrix, both of which are proportional to the genetic variance. If LD are estimated from the GWAS sample, $Var(\eta_j) = D_j^{-1} \left(\frac{m_j^0}{m} \sigma_g^2 + \sigma_e^2 \right)$. Further, if the genome-wide full LD matrix is used, $Var(\eta_j) = D_j^{-1} \sigma_e^2$, the same as that in (6).
3. If SNP j is independent of all other SNPs, $m_j^0 = m-1$ and $s_j^2 = 0$. Therefore, $Var(\eta_j) = D_j^{-1} \left(\frac{m-1}{m} \sigma_g^2 + \sigma_e^2 \right) = D_j^{-1} (\sigma_y^2 - \sigma_j^2) = D_j^{-1} (y'y - D_j \sigma_\beta^2) / n_j$, which is the residual variance under a single-SNP GWAS model.
4. Under some conditions, e.g., small \tilde{n}_j but large n_j , s_j^2 can be greater than 1. Thus, in the presence of LD sampling variance, the total residual variance (in the square brackets of (20)) can be greater than the phenotypic variance of the trait.

MCMC sampling scheme

The joint distribution of the data and parameters in model (7) is

$$\begin{aligned}
f(\mathbf{b}, \boldsymbol{\beta}, \pi, S, \sigma_\beta^2, \sigma_e^2) &\propto |\mathbf{R} \sigma_e^2|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma_e^2} (\mathbf{b} - \mathbf{W}\boldsymbol{\beta})' \mathbf{R}^{-1} (\mathbf{b} - \mathbf{W}\boldsymbol{\beta}) \right\} \\
&\times \prod_{j=1}^m \left[(h_j^S \sigma_\beta^2)^{-\frac{1}{2}} \exp \left\{ -\frac{\beta_j^2}{2h_j^S \sigma_\beta^2} \right\} \pi + \phi(1 - \pi) \right] \\
&\times \exp \left\{ -\frac{S^2}{2} \right\} \\
&\times (\sigma_\beta^2)^{-\frac{2+\nu_\beta}{2}} \exp \left\{ -\frac{\nu_\beta \tau_\beta^2}{2\sigma_\beta^2} \right\} \\
&\times (\sigma_e^2)^{-\frac{2+\nu_e}{2}} \exp \left\{ -\frac{\nu_e \tau_e^2}{2\sigma_e^2} \right\}
\end{aligned}$$

To obtain a joint posterior sample for parameter inference, we iteratively sample each parameter from its full conditional distribution. Except S , the full conditional distribution has a closed form for all the parameters, as shown below.

To deal with the mixture prior for β_j , we introduce an indicator variable δ_j

$$\delta_j \sim \text{Bernoulli}(\pi)$$

such that

$$\beta_j \begin{cases} \sim N(0, h_j^S \sigma_\beta^2), & \delta_j = 1 \\ = 0, & \delta_j = 0 \end{cases}$$

We first sample δ_j unconditional on β_j and then sample β_j conditional on δ_j , which has been shown to have slightly better mixing. The full conditional distribution for δ_j is

$$\delta_j | \mathbf{b}, \text{else} \sim \text{Bernoulli}(\hat{\pi})$$

where

$$\begin{aligned}\hat{\pi} &= \frac{f(\mathbf{b}|\delta_j = 1, else) \pi}{f(\mathbf{b}|\delta_j = 1, else) \pi + f(\mathbf{b}|\delta_j = 0, else) (1 - \pi)} \\ &= \left[1 + \frac{f(\mathbf{b}|\delta_j = 0, else) (1 - \pi)}{f(\mathbf{b}|\delta_j = 1, else) \pi} \right]^{-1}\end{aligned}\quad (21)$$

It is obvious that

$$f(\mathbf{b}|\delta_j = 0, else) = (2\pi |\mathbf{R}| \sigma_e^2)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma_e^2} \mathbf{b}'_{adj} \mathbf{R}^{-1} \mathbf{b}_{adj} \right\} \quad (22)$$

where $\mathbf{b}_{adj} = \mathbf{b} - \sum_{k \neq j} \mathbf{W}_k \beta_k$ is the adjusted \mathbf{b} for all the other SNP effects except SNP j . Fortunately, we do not need to compute this quantity because it will be cancelled in the likelihood ratio in (21), as shown below.

For $f(\mathbf{b}|\delta_j = 1, else)$, to be unconditional on β_j , we compute

$$\begin{aligned}f(\mathbf{b}|\delta_j = 1, else) &= \int f(\mathbf{b}|\delta_j = 1, \beta_j, else) f(\beta_j) d\beta_j \\ &= (2\pi |\mathbf{R}| \sigma_e^2)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma_e^2} \mathbf{b}'_{adj} \mathbf{R}^{-1} \mathbf{b}_{adj} \right\} \left(\frac{\sigma_{e_j}^{2*}}{h_j^S \sigma_\beta^2 C_j^*} \right)^{\frac{1}{2}} \exp \left\{ \frac{(r_j^*)^2}{2C_j^* \sigma_{e_j}^{2*}} \right\}\end{aligned}\quad (23)$$

This equation is derived based on integrating β_j out of the joint distribution of \mathbf{b} and β_j , which is closely related to the full conditional distribution of β_j , as shown next. Also see below for the definition of $\sigma_{e_j}^{2*}$, C_j^* and r_j^* . Substituting (22) and (23) into (21) gives the full conditional probability for $\delta_j = 1$.

Let $\sigma_{e_j}^{2*} = \left(\frac{n_j}{m} s_j^2 + \frac{m_j^0}{m} \right) \sigma_g^2 + \sigma_e^2$, which explicitly models the sampling variation of LD as described above.

The full conditional distribution of β_j (11) is

$$\beta_j | \mathbf{b}, \boldsymbol{\beta}_{-j}, else \sim N \left(\frac{r_j^*}{C_j^*}, \frac{\sigma_{e_j}^{2*}}{C_j^*} \right) \quad (24)$$

with r_j^* as in (13) (ignores t_j the cumulative effects of SNPs in chance LD) and $C_j^* = D_j + \frac{\sigma_{e_j}^{2*}}{h_j^S \sigma_\beta^2}$. It can be seen that instead of adjusting for the ‘‘leftover’’ effect from the mean, we shrink the mean towards zero while increase the variance (uncertainty) of the posterior distribution, because $\sigma_{e_j}^{2*}/C_j^* = 1 / \left(D_j / \sigma_{e_j}^{2*} + 1 / h_j^S \sigma_\beta^2 \right)$.

Given the sampled values of $\boldsymbol{\delta}$ and $\boldsymbol{\beta}$, the full conditional distribution for π , σ_β^2 and S is the same as in BayesS, with the sampling procedure elaborated in the Supplementary Note of Zeng et al [3].

The full conditional distribution for the residual variance σ_e^2 is

$$\sigma_e^2 | \mathbf{b}, else \sim \nu_e \tau_e^2 \chi_{\nu_e}^{-2} \quad (25)$$

where $\nu_e = \bar{n} + \nu_{e_0}$ and $\tau_e^2 = (\mathbf{e}'\mathbf{e} + \nu_{e_0} \tau_{e_0}^2) / \nu_e$ with ν_{e_0} and $\tau_{e_0}^2$ being the prior values. The residual sum of squares ($\mathbf{e}'\mathbf{e}$) can be computed as

$$\begin{aligned}
\mathbf{e}'\mathbf{e} &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\
&= \mathbf{y}'\mathbf{y} - 2\boldsymbol{\beta}'\mathbf{X}'\mathbf{y} + \boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta} \\
&= \mathbf{y}'\mathbf{y} - 2\boldsymbol{\beta}'\mathbf{r} + \boldsymbol{\beta}'(\mathbf{r} - \mathbf{r}_{adj}) \\
&= \mathbf{y}'\mathbf{y} - \boldsymbol{\beta}'\mathbf{r} - \boldsymbol{\beta}'\mathbf{r}_{adj}
\end{aligned} \tag{26}$$

where $\mathbf{r} = \mathbf{D}\mathbf{b}$ and $\mathbf{r}_{adj} = \mathbf{r} - \mathbf{X}'\mathbf{X}\boldsymbol{\beta}$ is the adjusted right-hand-side from the right-hand-side updating strategy (see below). The total sum of squares $\mathbf{y}'\mathbf{y}$ is computed from (9) and the median is used as the estimate of $\mathbf{y}'\mathbf{y}$ in (26).

The right-hand-side updating strategy and parallel computing

The MCMC implementation requires a computation of r_j in (12) for $m \times t$ times where m is the number of SNPs and t is the number of MCMC iterations. This is where the most majority computing time is spent. It can be seen from (12) that the summary-data level model is already much more efficient than the individual-data level model, because the vector-by-vector products $\mathbf{X}'_j\mathbf{y}$ and $\mathbf{X}'_j\mathbf{X}_k$ are replaced by scalar products $D_j b_j$ and $D_j^{\frac{1}{2}} B_{jk} D_k^{\frac{1}{2}}$. However, the adjustment of $D_j b_j$ (the right-hand-side \mathbf{r} of the mixed-model equations (10)) for all the other SNP effects is still too computationally intense given over a million of SNPs. To improve computational efficiency, we adopted the so-called right-hand-side updating algorithm for genomic prediction in the context of animal breeding [5]. We set out to compute a vector of adjusted right-hand-side

$$\mathbf{r}_{adj} = \mathbf{r} - \mathbf{C}\boldsymbol{\beta}$$

, where $\mathbf{r} = \mathbf{D}\mathbf{b}$ and $\mathbf{C} = \mathbf{D}^{\frac{1}{2}}\mathbf{B}\mathbf{D}^{\frac{1}{2}}$ (note that here \mathbf{C} is defined different from that in (10)). For each SNP, we compute

$$r_j = r_{adj,j} + C_{jj}\beta_j$$

and use r_j in the full conditional distribution of β_j (24). After a new value of β_j is sampled, we update the adjusted right-hand-side

$$\mathbf{r}_{adj}^{new} = \mathbf{r}_{adj}^{old} + \mathbf{C}_j(\beta_j^{old} - \beta_j^{new})$$

The benefit of this updating strategy comes from three sides. First, the computation of r_j for each SNP becomes trivial (reducing from vector to scalar operation). Second, the vector of \mathbf{r}_{adj} needs to be updated only when either β_j^{old} or β_j^{new} is not zero, thus the sparse genetic architecture will lead to a substantial gain in speed. Third, due to the use of a sparse LD matrix, the vector of \mathbf{C}_j has a large proportion of zero and therefore only a small fraction of \mathbf{r}_{adj} according to nonzero \mathbf{C}_j elements needs to be updated. In summary, the right-hand-side updating strategy substantially improves computational efficiency by taking the advantages of the sparse genetic architecture and the sparse LD correlation matrix. We further improve the efficiency by implementing a parallel

computing for sampling β_j of SNPs located on different chromosomes, and then combine results across threads to estimate the global parameters such as π , σ_β^2 , S and etc. This led to about 4 times faster when 4 cores were used with OpenMP library in our real trait analysis.

Algorithm pseudocode

Our SBayesS algorithm is implemented in GCTB based on the following pseudocode:

Algorithm 1 SBayesS algorithm

- 1: Initialise parameters and read summary statistics
 - 2: Reconstruct $\mathbf{X}'\mathbf{X}$ and $\mathbf{X}'\mathbf{y}$ from summary statistics and LD reference panel (\mathbf{D} = diagonal elements of $\mathbf{X}'\mathbf{X}$)
 - 3: Calculate $\mathbf{r}^* = \mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\boldsymbol{\beta}$
 - 4: **for** $i = 1$ to number of iterations **do**
 - 5: **for** $j = 1$ to number of SNPs **do**
 - 6: Calculate $\sigma_{e_j}^{2*} = \left(\frac{n_j}{m}s_j^2 + \frac{m_j^0}{m}\right)\sigma_g^2 + \sigma_e^2$
 - 7: Calculate the right hand side $r_j = r_j^* - D_j\beta_j$
 - 8: Calculate the left hand side $C_j = D_j + \frac{\sigma_{e_j}^{2*}}{(2p_jq_j)^S\sigma_\beta^2}$
 - 9: Calculate the posterior probability for $\delta_j = 1$ with $\Pr(\delta_j = 1|\mathbf{b}, \boldsymbol{\theta}) = \frac{1}{\exp[\log(L_0) - \log(L_1)]}$, where $\log(L_0) - \log(L_1) = \log(1 - \pi) - 0.5(\log(C_j^{-1}) - \log((2p_jq_j)^S\sigma_\beta^2) + C_j^{-1}r_j^2) - \log \pi$
 - 10: Sample δ_j based on the posterior probability
 - 11: **if** $\delta = 1$ **then**
 - 12: Sample β_j from the full conditional distribution $N\left(\frac{r_j}{C_j}, \frac{\sigma_{e_j}^{2*}}{C_j}\right)$
 - 13: Update $\mathbf{r}^* = \mathbf{r}^* + \mathbf{X}'\mathbf{X}_j(\beta_j^{(i-1)} - \beta_j)$, where $\mathbf{X}'\mathbf{X}_j$ is the j^{th} column of $\mathbf{X}'\mathbf{X}$ and $\beta_j^{(i-1)}$ is the sampled value from the last iteration
 - 14: **else**
 - 15: **if** $\beta_j^{(i-1)} = 0$ **then**
 - 16: Update $\mathbf{r}^* = \mathbf{r}^* + \mathbf{X}'\mathbf{X}_j\beta_j^{(i-1)}$
 - 17: **end if**
 - 18: **end if**
 - 19: **end for**
 - 20: Sample σ_β^2 from the full conditional distribution $\nu_\beta\tau_\beta^2\chi_{\nu_\beta}^{-2}$ with $\nu_\beta = m_{nz} + \nu_{\beta_0}$ and $\tau_\beta = \frac{\sum \frac{\beta_j^2}{(2p_jq_j)^S + \nu_{\beta_0}\tau_{\beta_0}}}{\nu_\beta}$
 - 21: Sample S from the full conditional distribution using HMC algorithm
 - 22: Sample π from the full conditional distribution $Beta(m_{NZ} + a_0, m - m_{NZ} + b_0)$
 - 23: Sample σ_e^2 from the full conditional distribution $\nu_e\tau_e^2\chi_{\nu_e}^{-2}$, where $\nu_e = \bar{n} + \nu_{e_0}$ and $\tau_e^2 = \frac{(\mathbf{e}'\mathbf{e} + \nu_{e_0}\tau_{e_0}^2)}{\nu_e}$ with $\mathbf{e}'\mathbf{e} = \mathbf{y}'\mathbf{y} - \boldsymbol{\beta}'\mathbf{X}'\mathbf{y} - \boldsymbol{\beta}'\mathbf{r}^*$
 - 24: Calculate the total genetic variance $\sigma_g^2 = \frac{\boldsymbol{\beta}'\mathbf{X}'\mathbf{y} - \boldsymbol{\beta}'\mathbf{r}^*}{\bar{n}}$
 - 25: Calculate the SNP-based heritability $h_{SNP}^2 = \frac{\sigma_g^2}{\sigma_g^2 + \sigma_e^2}$
 - 26: **end for**
-

Supplementary Note 2: Our genetic architecture parameter estimators are subject to SNP set and sample size

SNP-based heritability (h_{SNP}^2), polygenicity (π) and the relationship between MAF and effect size (S), all of which are defined with respect to a certain set of SNPs (see the definition of h_{SNP}^2 as an example [6]). Because the estimators for the genetic architecture parameters are based on the SNP markers rather than the actual causal variants, the estimates therefore rely on the tagging of the SNPs to the causal variants. For example,

our polygenicity estimate will be lower than that at the causal variants if many causal variants are not observed and poorly tagged by the SNPs (as shown in Supplementary Figure 4 in Zeng et al. [3]). In addition to the SNP panel, sample size can also affect our parameter estimates. As sample size increases, the polygenicity estimate is expected to increase because of the increased power to detect nonnull SNPs with very small effect sizes until all the genetic effects are detected. The other parameters, h_{SNP}^2 and S , are estimated based on the nonnull SNPs, therefore the estimators of these parameters are also SNP set and sample size dependent.

Despite the expected differences, the parameter estimates from the real traits were by-and-large consistent between different GWAS sample sizes. We randomly sampled 120k individuals from the UKBv3 dataset, performed GWAS and ran SBayesS. The correlation was 0.99 for SNP-heritability, 0.95 for polygenicity and 0.81 for the S parameter (Supplementary Fig. 6c). Compared to SNP-based heritability and polygenicity, the correlation for the S parameter was lower. In addition to the reason that S depends on the SNPs that are identified with nonzero effects and the number of which was larger with $N=350k$, the lower correlation for S was also because the S parameter is a hyper-parameter, which is more difficult to estimate and thereby estimated with a higher degree of uncertainty (i.e., a larger standard error).

To further investigate the property of our estimators, we performed a simulation based on the real trait results. That is, for each of the 18 traits in our benchmarking analysis, we simulated a trait using the estimated h_{SNP}^2 , π and S as the true values for the heritability, the proportion of causal variants and the S at the causal variants. Then, we ran SBayesS with the GWAS summary statistics computed with $n = 350K$ or $120K$, excluding the causal variants in the SNP panel. Supplementary Fig. 7a shows that \hat{h}_{SNP}^2 was slightly lower than the heritability at the causal variants (because of imperfect tagging), $\hat{\pi}$ is slightly higher than the proportion of causal variants (because of multiple SNPs jointly capturing a causal variant), and \hat{S} is almost unbiased. Supplementary Fig. 7b shows that the estimates with different sample sizes are highly concordant when the SNPs have a good tagging ability (in the simulation, we sampled causal variants from the 1.1 million common SNPs and left out the causal variants). This suggests that the additional variation and mean difference observed in the real data analysis (Supplementary Fig. 6c) are likely to be due to the insufficient tagging of HapMap3 common SNPs to the unobserved causal variants.

Finally, we ran SBayesS-strat with function annotations using the down-sampled dataset, and found a good consistency for different sample sizes (correlation $r = 0.99$ for h_{SNP}^2 enrichment, $r = 0.93$ for polygenicity enrichment, $r = 0.87$ for h_{NZE}^2 enrichment and $r = 0.73$ for the estimated S ; Supplementary Figure 31). Given a higher degree of uncertainty (a larger standard error), the S parameter estimates with the two different sample sizes were reasonably consistent, although the correlation for the S parameter was lower than that for h_{SNP}^2 and polygenicity as explained above. Compared to the per trait estimate (Supplementary Fig. 6c), the polygenicity enrichment estimates in the stratified model were much less sensitive to sample size because the extra number of small effects detected by a larger sample size in a functional category is likely proportional to that in the whole genome and therefore cancelled out in the ratio.

Supplementary Note 3: Simulation of GWAS summary statistics

The objective is to simulate the GWAS summary statistics $\widehat{\mathbf{b}}$ and their standard errors with an arbitrary sample size n given a reference sample. Consider model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$

Let \mathbf{D} be a diagonal matrix with $D_j = 2p_jq_jn_j$. The genotype matrix \mathbf{X} is centered but not standardised.

$$\mathbf{D}^{-1}\mathbf{X}'\mathbf{y} = \mathbf{D}^{-1}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} + \mathbf{D}^{-1}\mathbf{X}'\mathbf{e}$$

Let \mathbf{B} be the LD correlation matrix.

$$\begin{aligned}\widehat{\mathbf{b}} &= \mathbf{D}^{-\frac{1}{2}}\mathbf{B}\mathbf{D}^{\frac{1}{2}}\boldsymbol{\beta} + \mathbf{D}^{-1}\mathbf{X}'\mathbf{e} \\ &= \mathbf{b} + \boldsymbol{\epsilon}\end{aligned}$$

Assuming normality,

$$\widehat{\mathbf{b}} \sim N\left(\mathbf{D}^{-\frac{1}{2}}\mathbf{B}\mathbf{D}^{\frac{1}{2}}\boldsymbol{\beta}, \mathbf{D}^{-\frac{1}{2}}\mathbf{B}\mathbf{D}^{-\frac{1}{2}}\sigma_\epsilon^2\right)$$

The LD matrix \mathbf{B} is estimated from the reference sample of a relatively small sample size. In order to remove chance LD and facilitate computation, we follow the method in O'Connor and Price [7] to convert the full matrix into a positive semidefinite block diagonal matrix. The first step is to convert the full matrix into a block diagonal matrix with 50 blocks, each containing a LD window size of 2Mb (the total chromosome length is 100Mb in the simulation). The next step is to perform an eigen decomposition on each block:

$$\mathbf{B}_i = \mathbf{V}_i\boldsymbol{\Sigma}_i\mathbf{V}_i'$$

To remove noise and ensure a positive definite LD matrix, we removed the non-positive eigenvalues in $\boldsymbol{\Sigma}_i$ and the corresponding eigenvectors in \mathbf{V}_i . The last step is to re-normalize each block to have the diagonal values equal to one:

$$\mathbf{B}_i^* = \mathbf{J}_i^{-\frac{1}{2}}\mathbf{B}_i\mathbf{J}_i^{-\frac{1}{2}}$$

where \mathbf{J}_i is the diagonal matrix corresponding to the diagonal of \mathbf{B}_i .

Thus,

$$\mathbf{b} = \mathbf{D}^{-\frac{1}{2}}\mathbf{B}^*\mathbf{D}^{\frac{1}{2}}\boldsymbol{\beta}$$

Since

$$\mathbf{B}_i^* = \mathbf{J}_i^{-\frac{1}{2}}\mathbf{V}_i\boldsymbol{\Sigma}_i^{\frac{1}{2}}\boldsymbol{\Sigma}_i^{\frac{1}{2}}\mathbf{V}_i'\mathbf{J}_i^{-\frac{1}{2}}$$

let $u \sim N(0, 1)$, then

$$\text{Var} \left(\mathbf{D}_i^{-\frac{1}{2}} \mathbf{J}_i^{-\frac{1}{2}} \mathbf{V}_i \boldsymbol{\Sigma}_i^{\frac{1}{2}} \mathbf{u}_i \sigma_\epsilon \right) = \mathbf{D}_i^{-\frac{1}{2}} \mathbf{B}_i^* \mathbf{D}_i^{-\frac{1}{2}} \sigma_\epsilon^2$$

Thus,

$$\boldsymbol{\epsilon}_i = \mathbf{D}_i^{-\frac{1}{2}} \mathbf{J}_i^{-\frac{1}{2}} \mathbf{V}_i \boldsymbol{\Sigma}_i^{\frac{1}{2}} \mathbf{u}_i \sigma_\epsilon$$

where $\sigma_\epsilon^2 \approx \sigma_y^2$. The standard error of \widehat{b}_j (SE) is the square root of

$$\text{Var} \left(\widehat{b}_j \right) = D_j^{-1} \sigma_{\epsilon,j}^2 = D_j^{-1} \left(\sigma_y^2 - 2p_j q_j \widehat{b}_j^2 \right)$$

Supplementary Note 4: Maximum likelihood estimation of S

Here, we derive the maximum likelihood estimate (MLE) of S given the true effect sizes ($\boldsymbol{\beta}$) of m causal variants.

Since our model is

$$\beta_j \sim N \left(0, h_j^S \sigma_\beta^2 \right)$$

where $h_j = 2p_j(1 - p_j)$, the likelihood function is

$$f \left(\boldsymbol{\beta} | S, \sigma_\beta^2 \right) = \prod_{j=1}^m (2\pi)^{-\frac{1}{2}} \left(h_j^S \sigma_\beta^2 \right)^{-\frac{1}{2}} \exp \left\{ -\frac{\beta_j^2}{2h_j^S \sigma_\beta^2} \right\}$$

Thus, the log-likelihood function, after dropping out the normalising constant, is

$$\mathcal{L} = -\frac{S}{2} \sum_{j=1}^m \log h_j - \frac{m}{2} \log \sigma_\beta^2 - \frac{1}{2\sigma_\beta^2} \sum_{j=1}^m \frac{\beta_j^2}{h_j^S}$$

Taking the partial derivatives of \mathcal{L} with respect to S and σ_β^2 , respectively, gives

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial S} &= -\frac{1}{2} \sum_{j=1}^m \log h_j + \frac{1}{2\sigma_\beta^2} \sum_{j=1}^m \left(\frac{\beta_j^2}{h_j^S} \log h_j \right) \\ \frac{\partial \mathcal{L}}{\partial \sigma_\beta^2} &= -\frac{m}{2\sigma_\beta^2} + \frac{1}{2(\sigma_\beta^2)^2} \sum_{j=1}^m \frac{\beta_j^2}{h_j^S} \end{aligned}$$

The MLE of S and σ_β^2 can be obtained by solving the above equations that set to be zero. In the simulation study, we used *optim* function in R to find the MLE of S .

Supplementary Note 5: Standard error of predicted evolutionary parameter

We estimated the standard error of predicted evolutionary parameter by considering a random coefficient linear model. Our model for predicting an evolutionary parameter is: $\widehat{\theta} = \widehat{\mathbf{X}}' \widehat{\boldsymbol{\beta}}$ where $\widehat{\mathbf{X}}$ is a vector of the three genetic architecture parameter estimates from SBayesS and their respective polynomial terms and $\widehat{\boldsymbol{\beta}}$ is a vector of the

weights estimated from the reference dataset of simulation. There are two sources of estimation uncertainty. The first one is the uncertainty in estimating $\widehat{\mathbf{X}}$ by SBayesS. The second source is the uncertainty in estimating $\widehat{\boldsymbol{\beta}}$ by the polynomial regression in training. The estimation variance of $\widehat{\boldsymbol{\theta}}$ can be therefore written as

$$\begin{aligned}
\text{Var}(\widehat{\mathbf{X}}'\widehat{\boldsymbol{\beta}}) &= \text{Var}_{\widehat{\boldsymbol{\beta}}}\left[E_{\widehat{\mathbf{X}}}\left(\widehat{\mathbf{X}}'\widehat{\boldsymbol{\beta}}\mid\widehat{\boldsymbol{\beta}}\right)\right] + E_{\widehat{\boldsymbol{\beta}}}\left[\text{Var}_{\widehat{\mathbf{X}}}\left(\widehat{\mathbf{X}}'\widehat{\boldsymbol{\beta}}\mid\widehat{\boldsymbol{\beta}}\right)\right] \\
&= \text{Var}_{\widehat{\boldsymbol{\beta}}}\left[\mathbf{X}'\widehat{\boldsymbol{\beta}}\right] + E_{\widehat{\boldsymbol{\beta}}}\left[\widehat{\boldsymbol{\beta}}'\text{Var}\left(\widehat{\mathbf{X}}\right)\widehat{\boldsymbol{\beta}}\right] \\
&= \mathbf{X}'\text{Var}\left(\widehat{\boldsymbol{\beta}}\right)\mathbf{X} + E_{\widehat{\boldsymbol{\beta}}}\left[\text{tr}\left(\widehat{\boldsymbol{\beta}}'\text{Var}\left(\widehat{\mathbf{X}}\right)\widehat{\boldsymbol{\beta}}\right)\right] \\
&= \mathbf{X}'\text{Var}\left(\widehat{\boldsymbol{\beta}}\right)\mathbf{X} + E_{\widehat{\boldsymbol{\beta}}}\left[\text{tr}\left(\widehat{\boldsymbol{\beta}}\widehat{\boldsymbol{\beta}}'\text{Var}\left(\widehat{\mathbf{X}}\right)\right)\right] \\
&= \mathbf{X}'\text{Var}\left(\widehat{\boldsymbol{\beta}}\right)\mathbf{X} + \text{tr}\left[E\left(\widehat{\boldsymbol{\beta}}\widehat{\boldsymbol{\beta}}'\right)\text{Var}\left(\widehat{\mathbf{X}}\right)\right] \\
&= \mathbf{X}'\text{Var}\left(\widehat{\boldsymbol{\beta}}\right)\mathbf{X} + \text{tr}\left[\text{Var}\left(\widehat{\boldsymbol{\beta}}\right)\text{Var}\left(\widehat{\mathbf{X}}\right)\right]
\end{aligned}$$

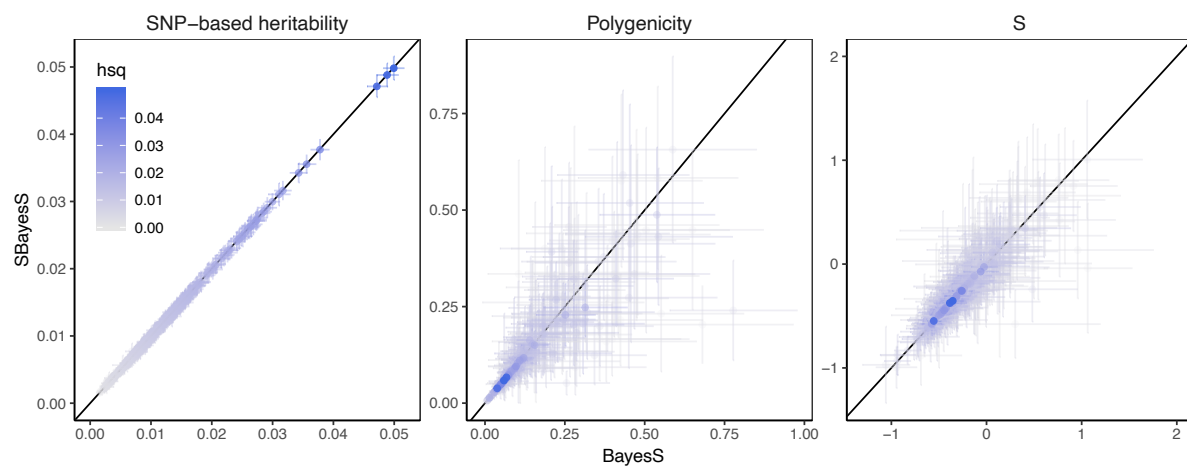
where $\text{Var}\left(\widehat{\mathbf{X}}\right)$ is the variance-covariance matrix of the posterior estimates computed from the MCMC samples of SBayesS and $\text{Var}\left(\widehat{\boldsymbol{\beta}}\right)$ is the estimation variance matrix of the weights obtained from the polynomial regression in the reference dataset. The true values of \mathbf{X} are unknown and therefore we replace them by their estimates $\widehat{\mathbf{X}}$.

Supplementary Note 6: Acknowledgements

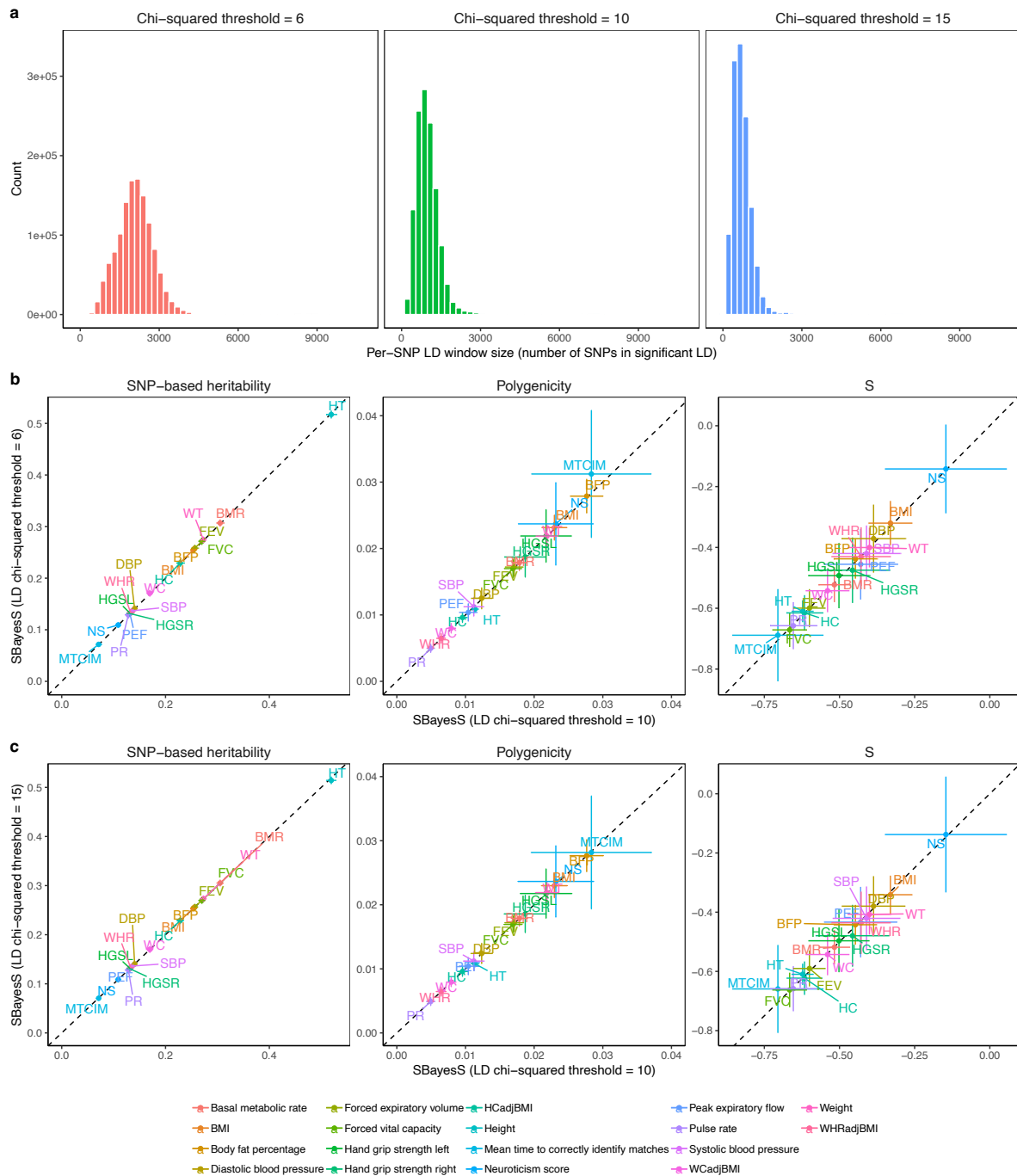
UKB: This study has been conducted using UK Biobank resource under Application Number 12505. UK Biobank was established by the Wellcome Trust medical charity, Medical Research Council, Department of Health, Scottish Government and the Northwest Regional Development Agency. It has also had funding from the Welsh Assembly Government, British Heart Foundation and Diabetes UK.

GERA: The Genetic Epidemiology Research on Adult Health and Aging study was supported by grant RC2 AG036607 from the National Institutes of Health, grants from the Robert Wood Johnson Foundation, the Ellison Medical Foundation, the Wayne and Gladys Valley Foundation and Kaiser Permanente. The authors thank the Kaiser Permanente Medical Care Plan, Northern California Region (KPNC) members who have generously agreed to participate in the Kaiser Permanente Research Program on Genes, Environment and Health (RPGEH).

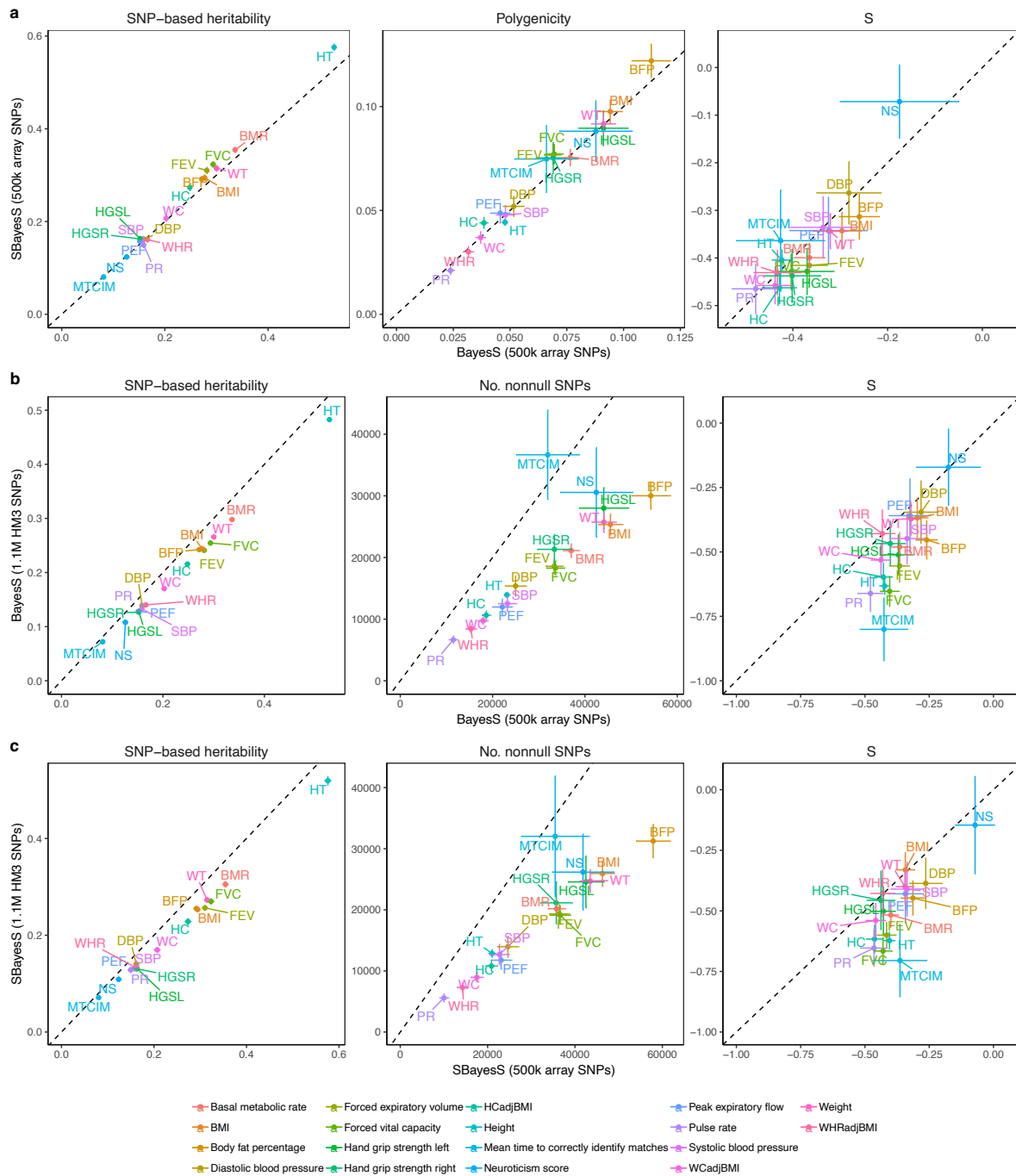
Supplementary Figures



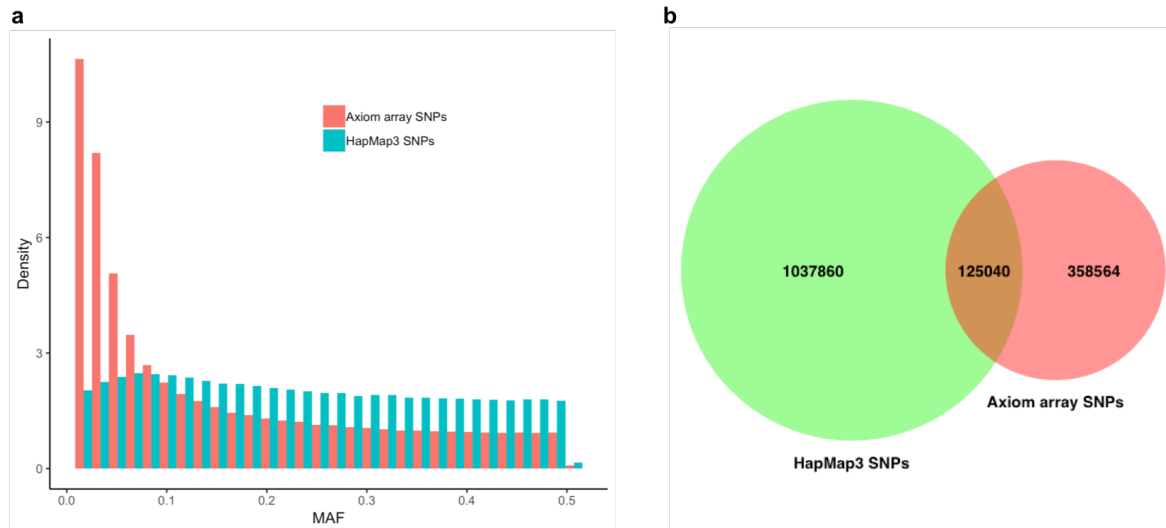
Supplementary Figure 1 Benchmarking SBayesS with BayesS using the same data in the chromosome-wide analysis across 18 UKB traits. The comparison was based on the unrelated individuals of European ancestry in the interim release of the UKB data (max n=120k) and ~500k array genotyped common SNPs (MAF>0.01). In the SBayesS analysis, the full LD matrix that included all pairwise LD was used for each chromosome. Data are presented as posterior means +/- posterior standard errors.



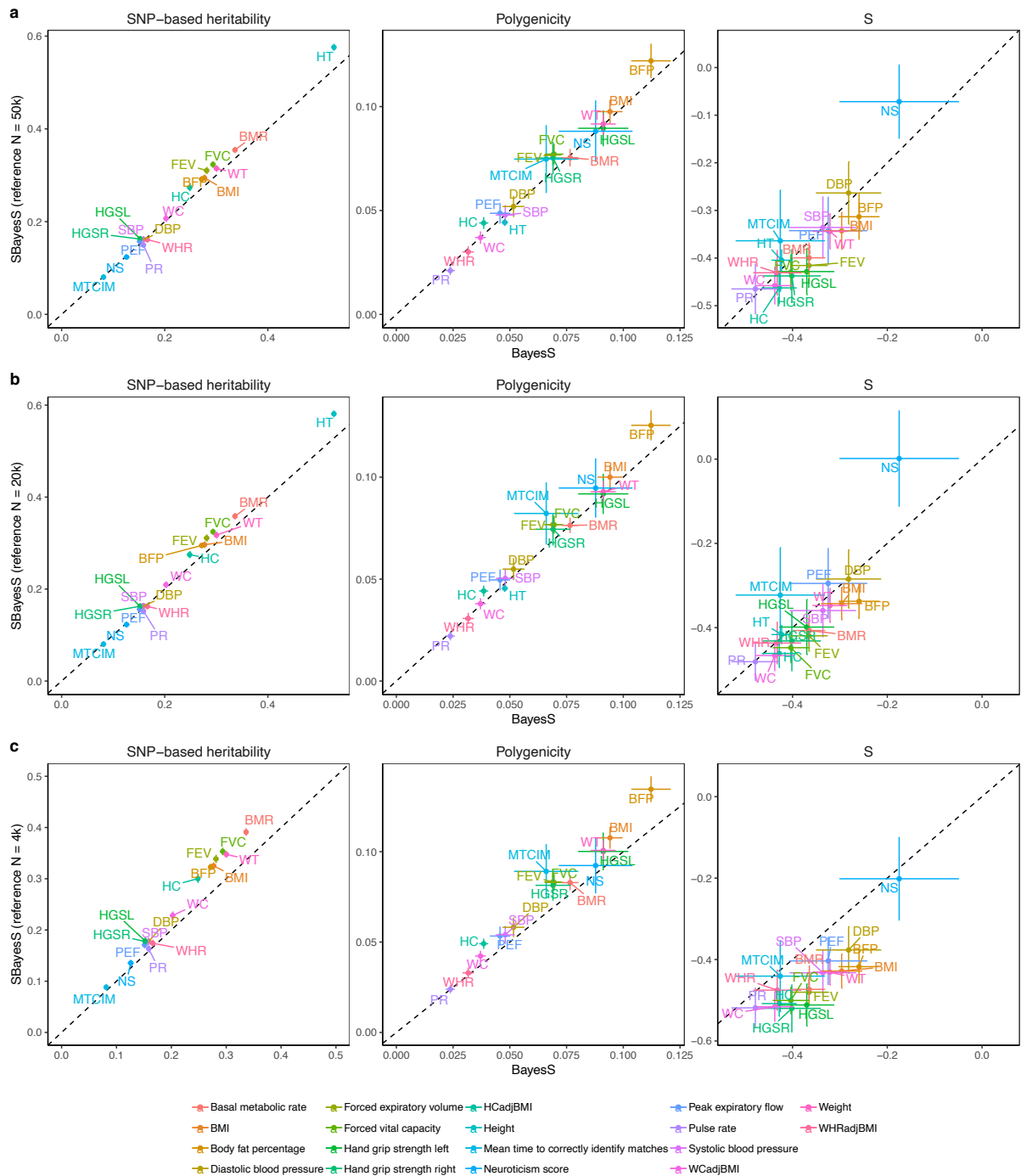
Supplementary Figure 2 Assessing the performance of SBayesS with different chi-squared thresholds used to make the sparse LD matrix. We computed summary statistics for 18 traits using the interim release of the UKB data (max $n=120k$) with ~ 1.1 million HapMap3 common SNPs. The sparse LD matrix was computed from a random sample of 50k unrelated individuals in the full UKB data with a chi-squared threshold of 6, 10 or 15 (corresponding to a r^2 threshold of 1, 2 or 3×10^{-4} , respectively). a) Distributions of the numbers of SNPs detected in LD with the target SNP given different chi-squared thresholds; b) Comparison between SBayesS results with chi-squared thresholds of 6 and 10; c) Comparison between SBayesS results with chi-squared thresholds of 10 and 15. Data are presented as posterior means \pm posterior standard errors.



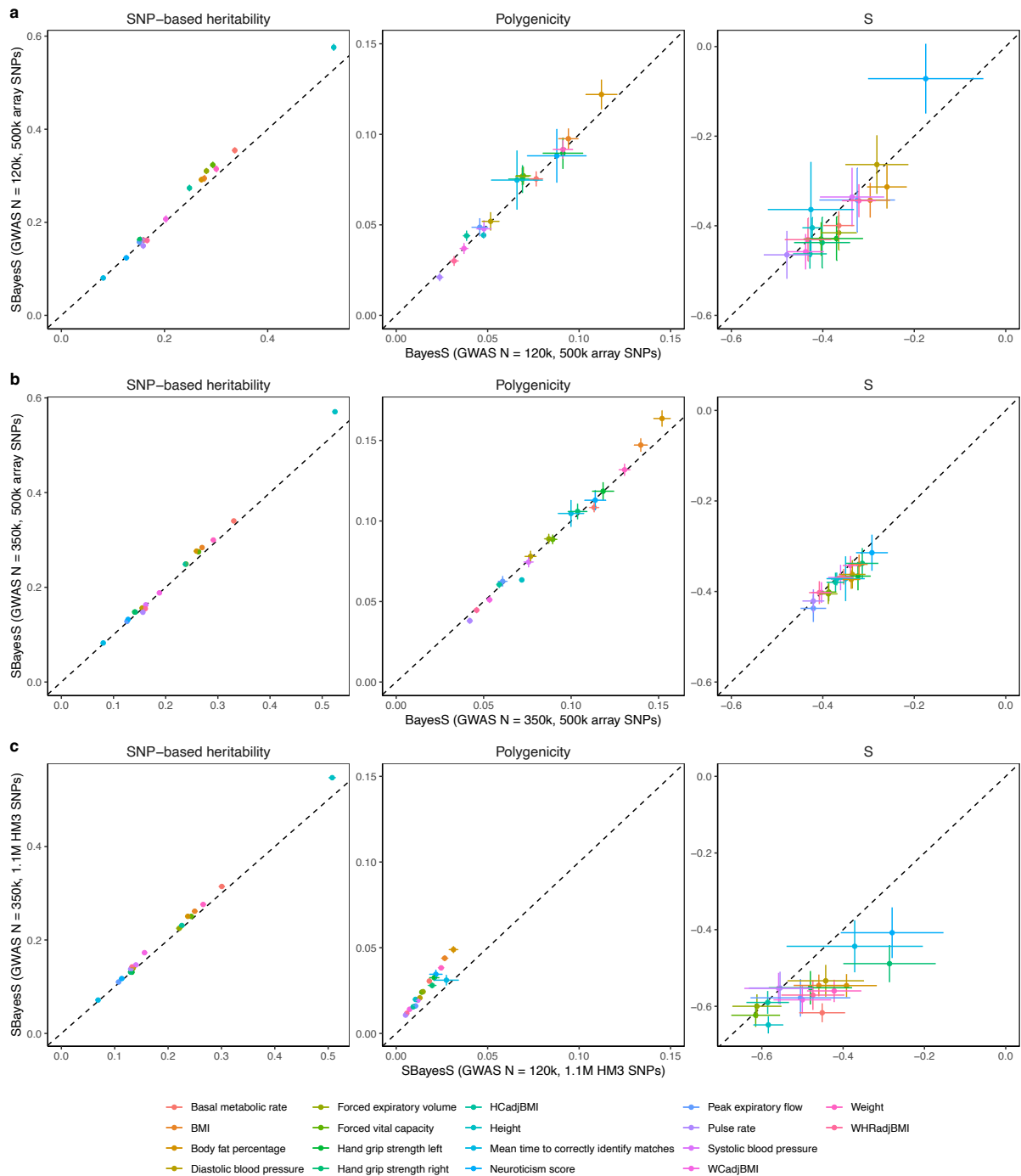
Supplementary Figure 3 Benchmarking SBayesS with BayesS given different SNP panels. We used the unrelated individuals of European ancestry in the interim release of the UKB data (max $n=120k$) and two SNP panels ($\sim 500k$ Affymetrix array SNPs and ~ 1.1 million HapMap3 SNPs) for the SBayesS analysis. The sparse LD matrix was computed from a random sample of 50k unrelated individuals from the full UKB cohort at a chi-squared threshold of 10. a) Comparison between SBayesS and BayesS using array SNPs; b) Comparison between BayesS results using HapMap3 and array SNPs; c) Comparison between SBayesS results using HapMap3 and array SNPs. For a fair comparison of π between panels, the number of SNPs with nonzero effects is shown in b) and c). Data are presented as posterior means \pm posterior standard errors.



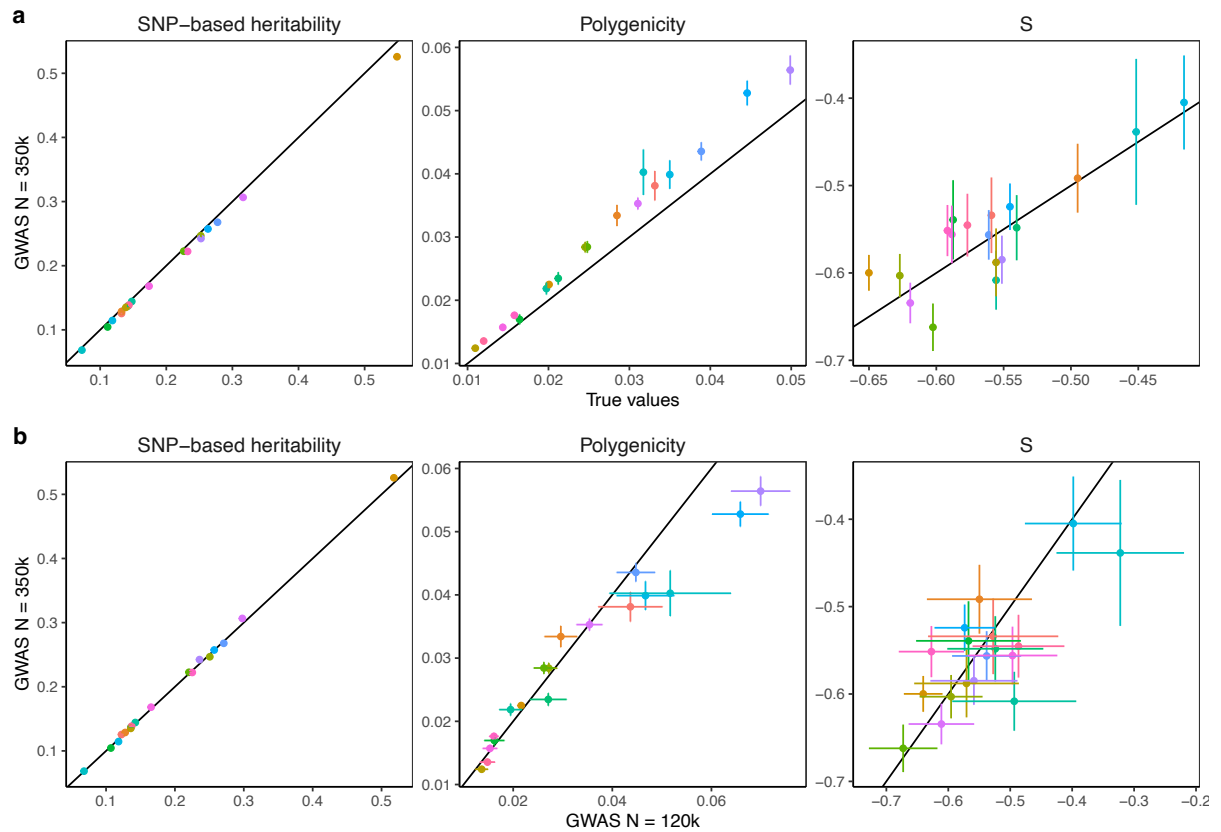
Supplementary Figure 4 Distributions of MAF of common SNPs (MAF>1%) in Affymetrix Axiom array and HapMap3 from the interim release of the UKB data. The array SNP panel was more enriched with low-frequency SNPs and had only a small overlap with the HapMap3 SNP panel. This may explain the differences in genetic architecture parameter estimates between array and HapMap3 SNPs (Supplementary Fig. 3). For example, the array SNPs might be more efficient to capture the low-frequency variance and therefore had slightly higher total SNP-based heritability. However, the majority of the common SNPs in HapMap3 Panel is believed to have better tagging to the causal variants. Thus, it is reasonable that the polygenicity estimates were lower than those with array SNPs because the model does not need multiple SNPs in low LD with the causal variants to jointly capture the causal effects. Similarly, a stronger estimate of S is expected because if the causal effects spread on multiple SNPs in LD, each SNP would have a relatively small effect size, which will dilute the signal for the relationship between effect size and MAF.



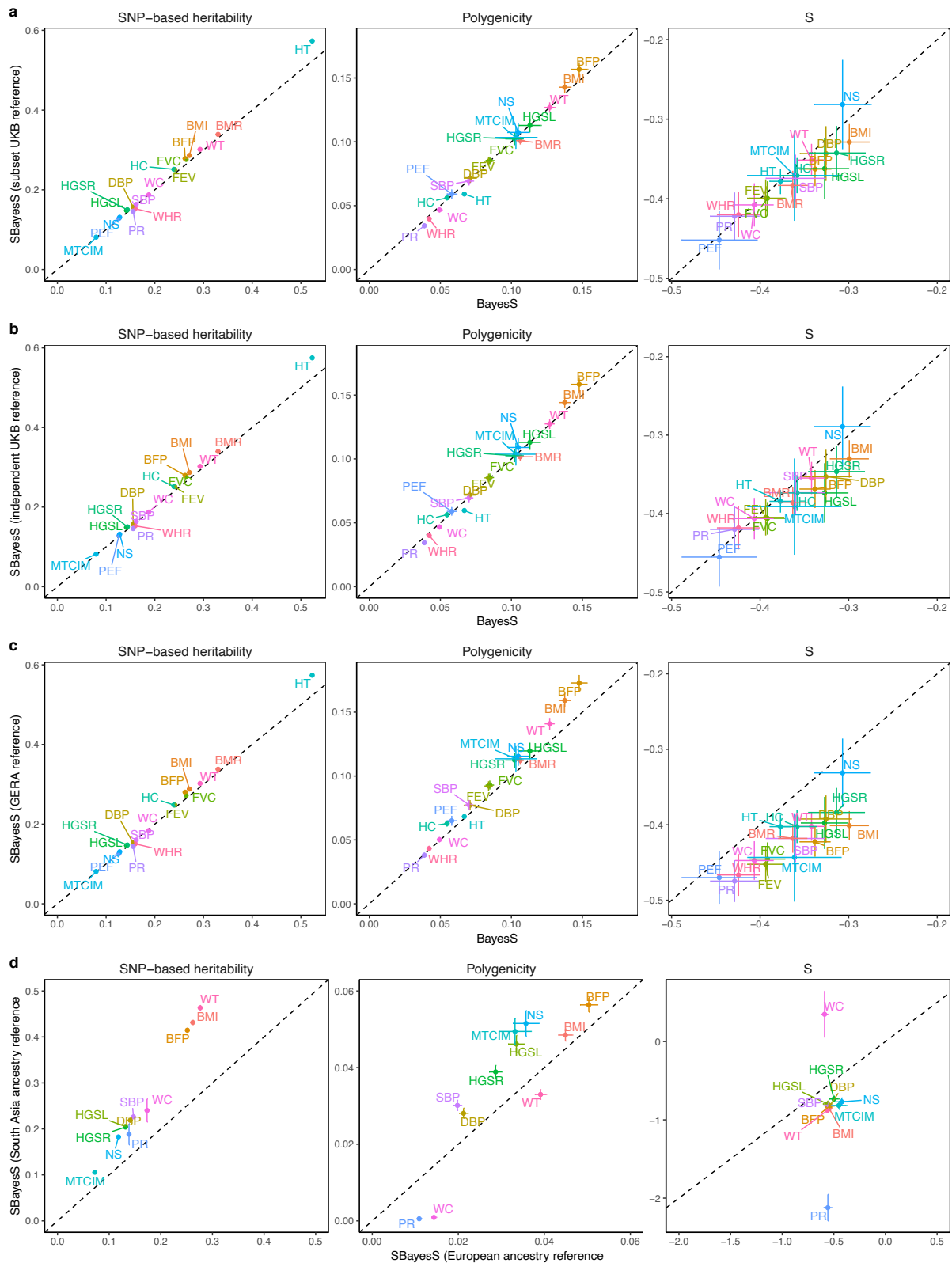
Supplementary Figure 5 Benchmarking SBayesS with BayesS given different reference sample sizes. We used the unrelated individuals of European ancestry in the interim release of the UKB data (max n=120k) and ~500k Affymetrix array common SNPs. The sparse LD matrix was computed from a random sample of a) 50k, b) 20k or c) 4k unrelated individuals from the full UKB data at a chi-square threshold of 10. The inflation in parameter estimation increased when the reference sample size was too small. When the reference sample size was 4k, SBayesS analysis for height did not converge. Data are presented as posterior means +/- posterior standard errors.



Supplementary Figure 6 Benchmarking SBayesS with BayesS given different GWAS sample sizes. a) and b) Comparison between SBayesS and BayesS using the unrelated individuals of European ancestry in the interim (max n=120k) and full (max n=350k) releases of the UKB data and ~500k Affymetrix array SNPs. c) Comparison between SBayesS results given GWAS sample size of 350k and 120k using ~1.1 million HapMap3 SNPs. The sparse LD matrix used in SBayesS was computed from a random sample of 50k unrelated individuals from the full UKB data at a chi-squared threshold of 10. Data are presented as posterior means +/- posterior standard errors.



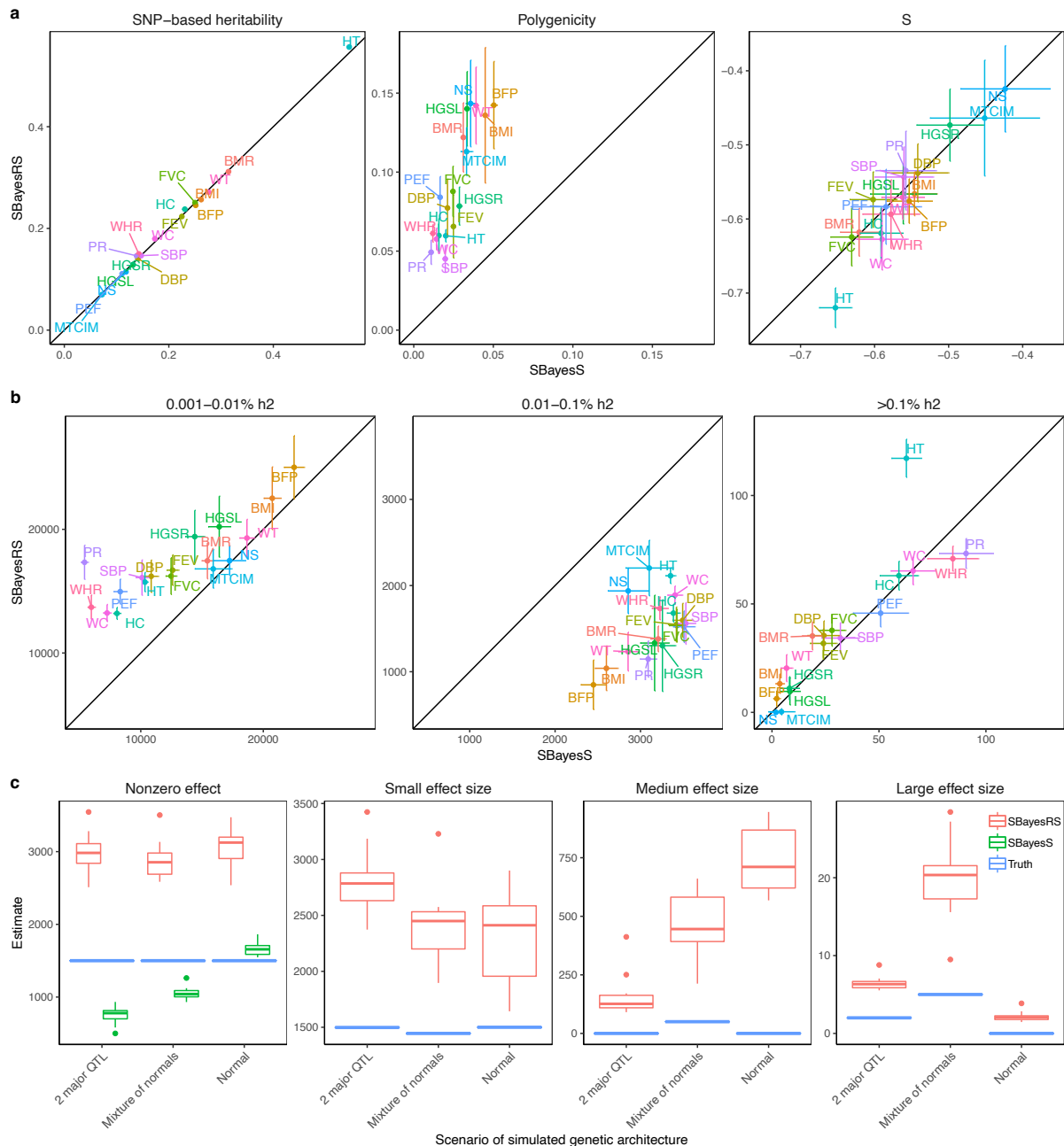
Supplementary Figure 7 SBayesS results of the simulation study based on the parameter estimates from the real trait analysis. a) Comparison of parameter estimates based on the SNP markers (causal variants excluded) and the true parameter values at the causal variants in the simulation. b) Comparison of parameter estimates with GWAS sample sizes of 350K and 120K. Data are presented as posterior means +/- posterior standard errors.



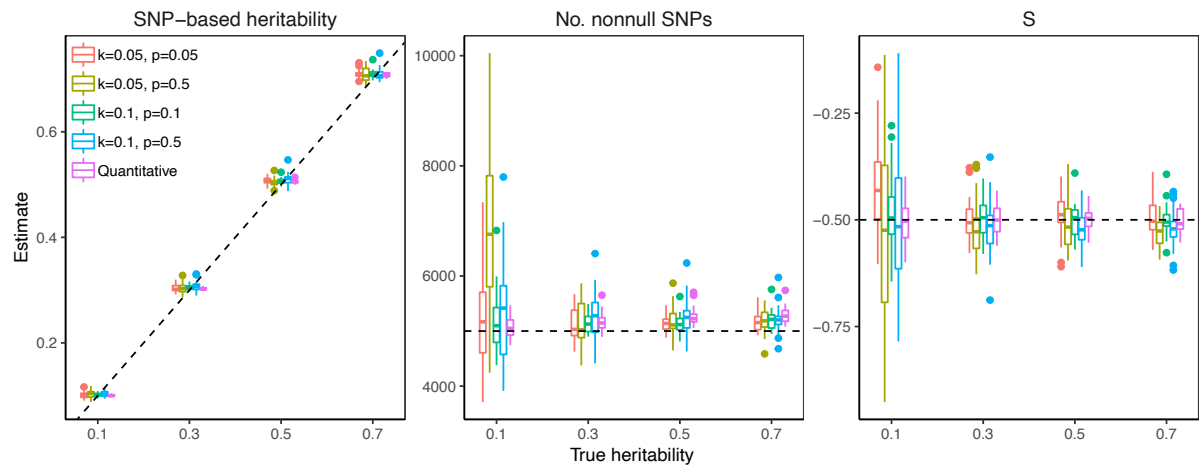
Supplementary Figure 8 Benchmarking SBayesS with BayesS given different LD references.

We used phenotypes from a random sample of 300k unrelated individuals of European ancestry from the full UKB data and ~500k Affymetrix array SNPs. The sparse LD matrix with a chi-squared threshold of 10 used in SBayesS was computed from a) a subset sample of 50k UKB individuals, b) an independent sample of 50k UKB individuals, or c) 50k unrelated individuals

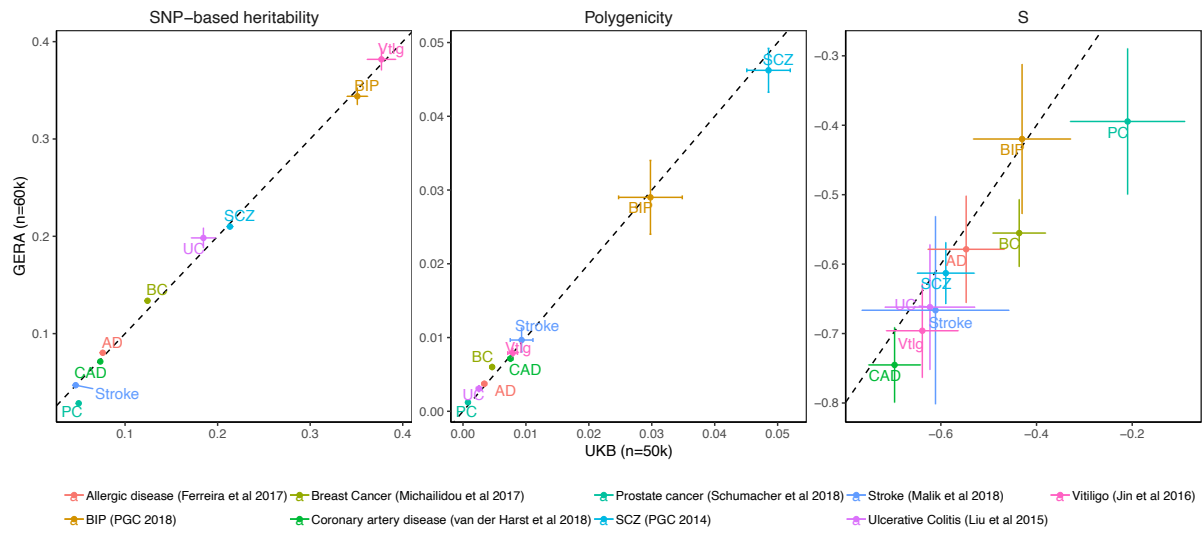
from the GERA dataset. d) When GWAS data were from the UKB European population, using a sample ($n = 9948$) of South Asia ancestry as LD reference in SBayesS (1.1 million HapMap3 SNPs) resulted in a severe bias to the genetic architecture parameter estimates and a failure in convergence for 11 out of the 18 traits (61%) with a chromosome-wide full LD matrix (all traits were failed in convergence when using a sparse LD matrix with the default sparsity, i.e., chi-squared statistic threshold of 10). Data are presented as posterior means \pm posterior standard errors.



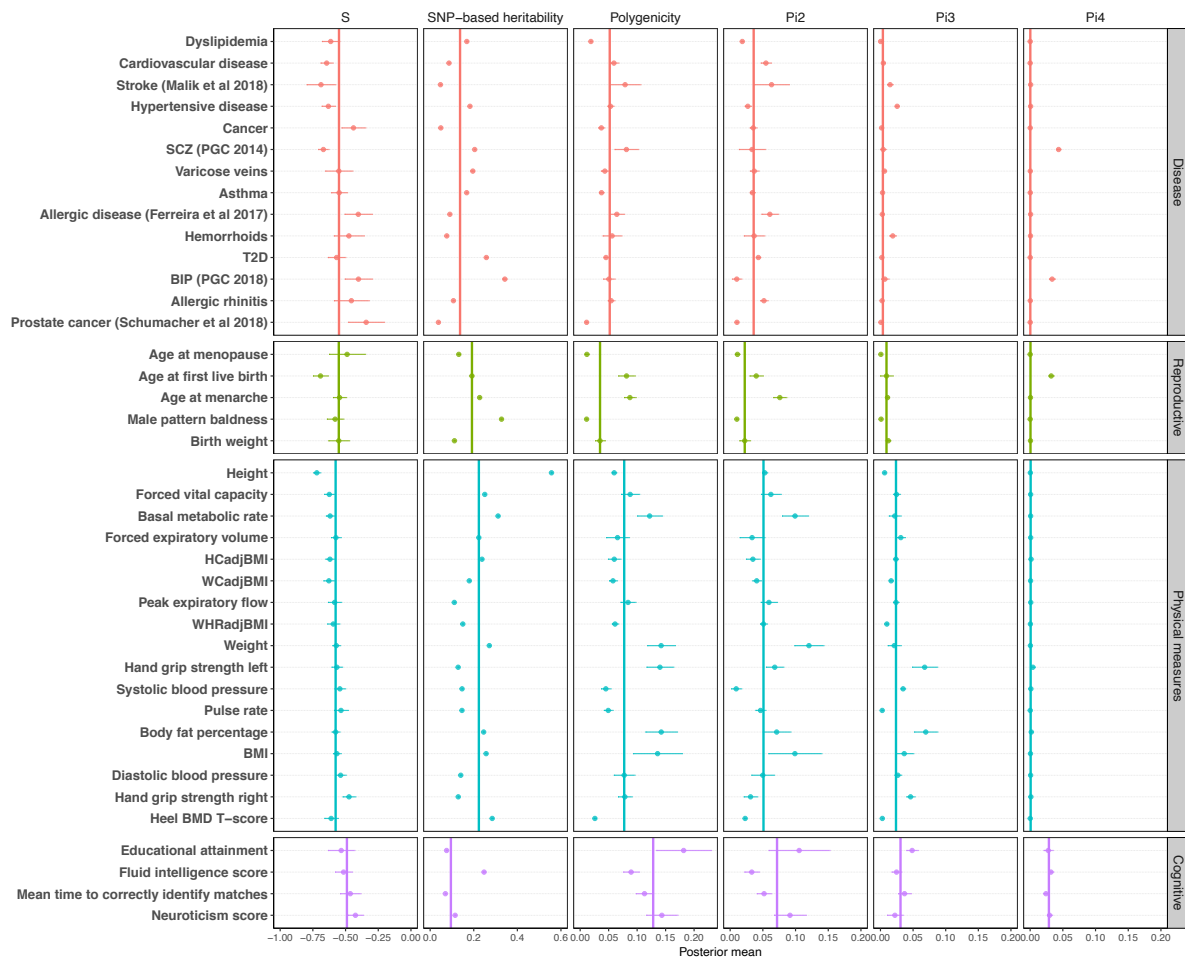
Supplementary Figure 9 Comparison between SBayesS and SBayesRS for the analyses of 18 quantitative traits in the UKB. a) Estimates of genetic architecture parameters, i.e. S , SNP-based heritability and polygenicity. b) Estimated number of SNPs that explain 0.001-0.01%, 0.01-0.1% and >0.1% of the total SNP-based heritability in SBayesRS and SBayesS. Data are presented as posterior means +/- posterior standard errors. c) Estimated number of SNPs with nonzero effects as well as that in the small, medium and large mixture components in the simulation on chromosomes 21 and 22. The band inside the box is the median, the bottom and top of the box are the first and third quartiles, respectively (Q1 and Q3), and the lower and upper whiskers are $Q1 - 1.5 \text{ IQR}$ and $Q3 + 1.5 \text{ IQR}$, respectively, where $\text{IQR} = Q3 - Q1$.



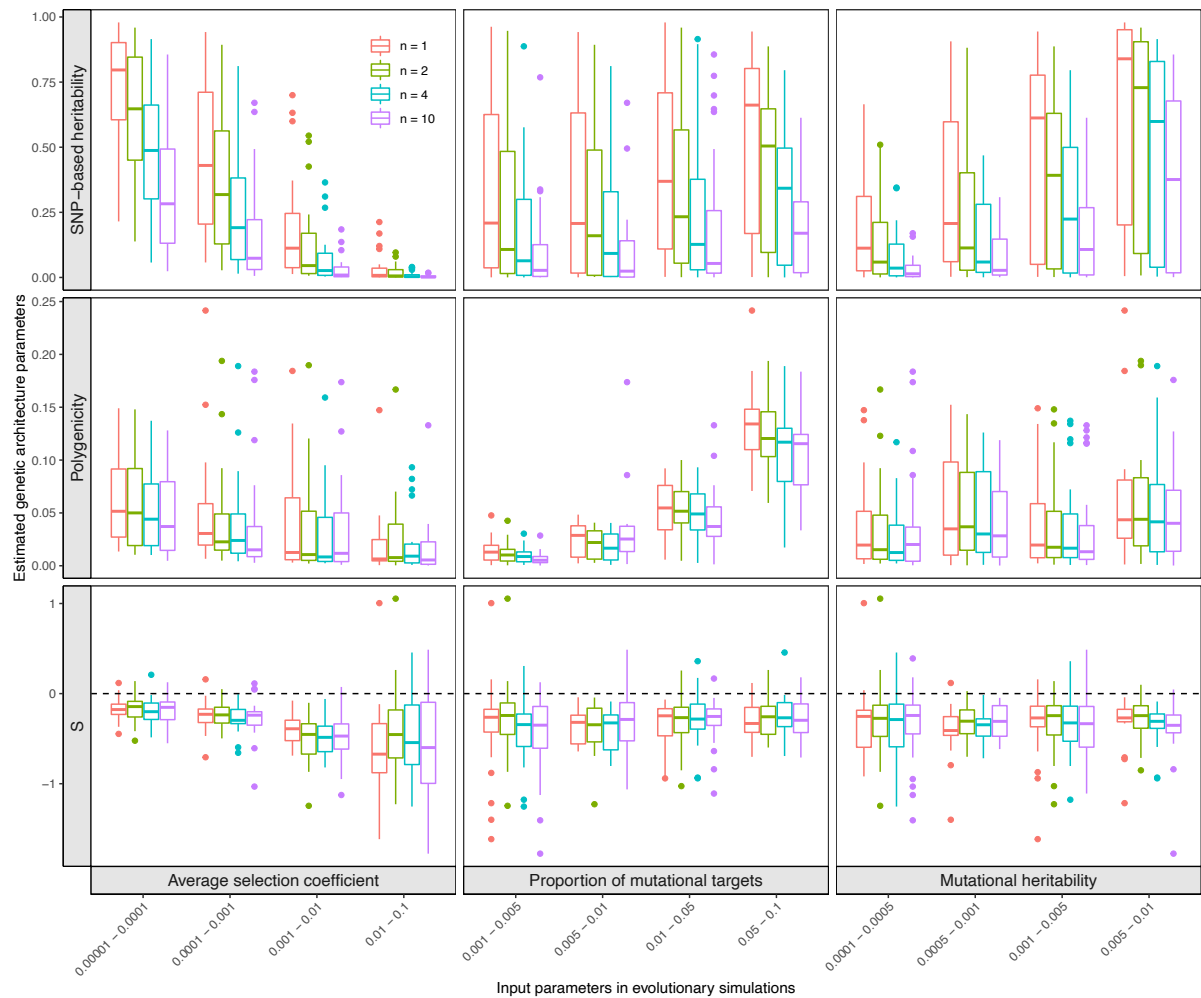
Supplementary Figure 10 Estimation of the three genetic architecture parameters using SBayesS with simulated data for quantitative traits or case-control studies with different population (k) and sample (p) prevalence. The band inside the box is the median, the bottom and top of the box are the first and third quartiles, respectively (Q1 and Q3), and the lower and upper whiskers are $Q1 - 1.5 \text{ IQR}$ and $Q3 + 1.5 \text{ IQR}$, respectively, where $\text{IQR} = Q3 - Q1$. We used the full UKB data ($n=350k$) for simulations, where 5k SNPs were randomly chosen from ~ 1.1 million HapMap3 common SNPs as causal variants with true $S=-0.5$, and the trait heritability was set to be 0.1, 0.3, 0.5 or 0.7 (at the liability scale for the binary trait). The sparse LD matrix used in SBayesS was computed from a subset sample of 50k UKB individuals with a chi-squared threshold of 10. The simulation was repeated 30 times for each scenario. When $k=0.05$, $p=0.5$ and true heritability=0.1, the polygenicity estimate tend to bias upward with large estimation variation, which is likely due to insufficient power to distinguish the model that fits only causal variants from that fits multiple SNPs in LD with the causal variants.



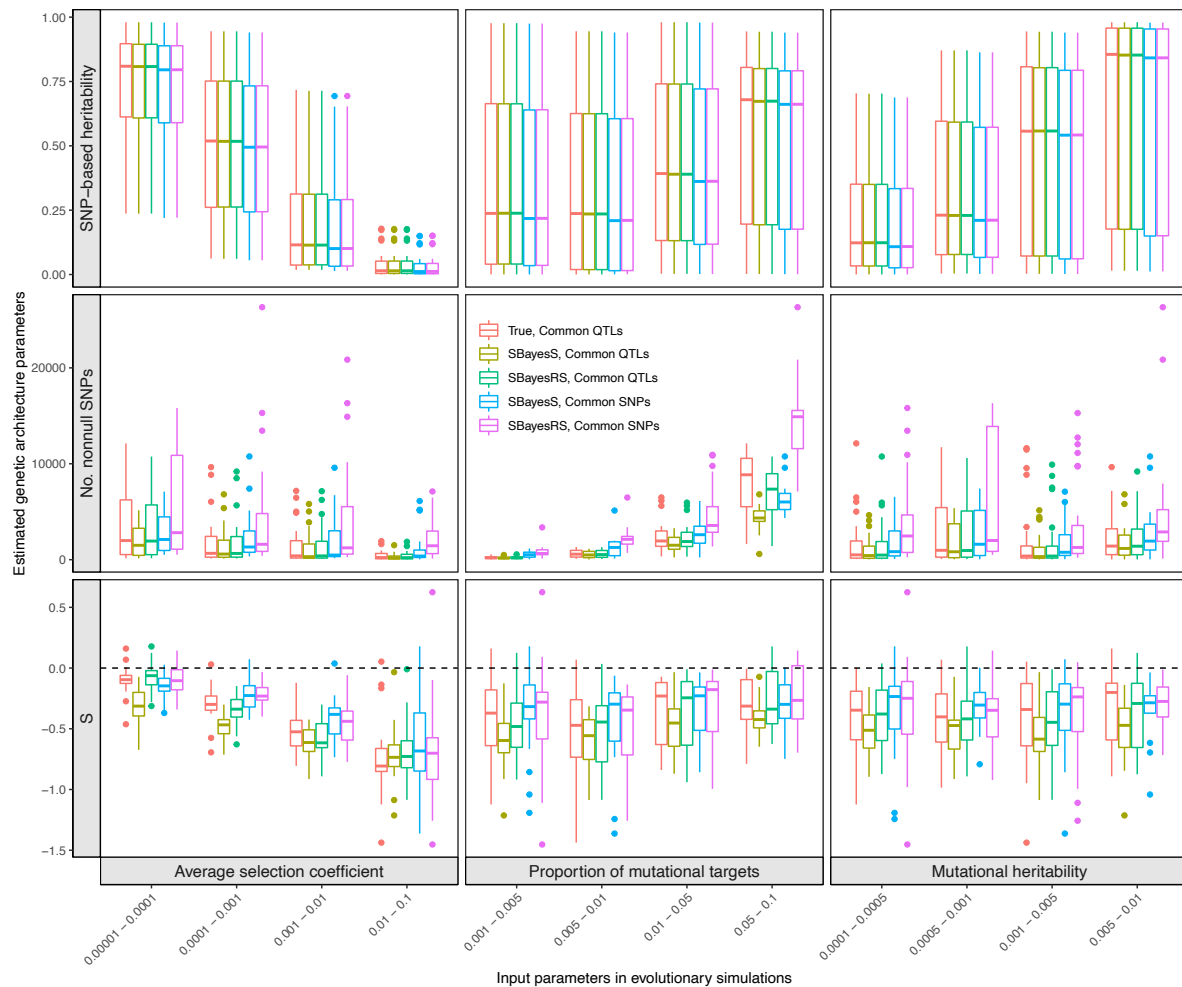
Supplementary Figure 12 Comparison between SBayesS results based on LD from GERA and UKB (a random subsample of 50k unrelated individuals) using published GWAS summary data for 9 diseases. Data are presented as posterior means +/- posterior standard errors. Colours with acronyms indicate different traits, whose full names are shown at the bottom of the figure.



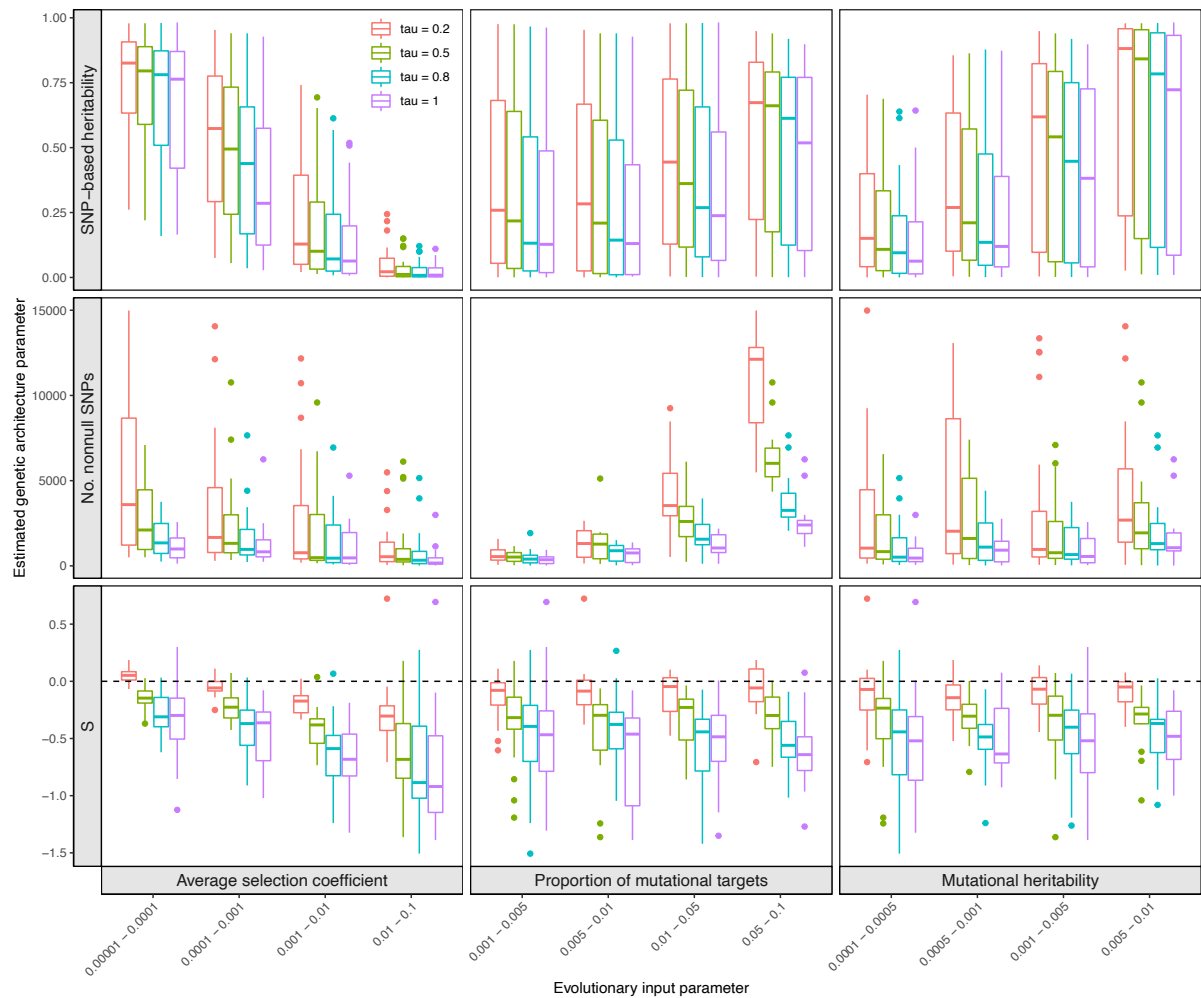
Supplementary Figure 13 Estimation of the three genetic architecture parameters for 35 complex traits (including diseases) in UKB (max n=350k) and 5 common diseases from published GWAS (labelled with publications) by SBayesRS. Shown are the posterior means (dots) and standard errors (horizontal bars) of the parameters for each trait. The colour indicates the category that the trait belongs to. The vertical bar shows the median of the estimates across traits in each category. Pi2, Pi3 and Pi4 show the proportions of SNPs in the small, medium and large effect size components of the mixture distribution in SBayesRS. Traits are in the same order as in Figure 2 for comparison. Coronary artery disease (van der Harst et al 2018), Vitiligo (Jin et al 2016) and Ulcerative Colitis (Liu et al 2015) and Breast cancer (Michallidou et al 2017) did not have converged results.



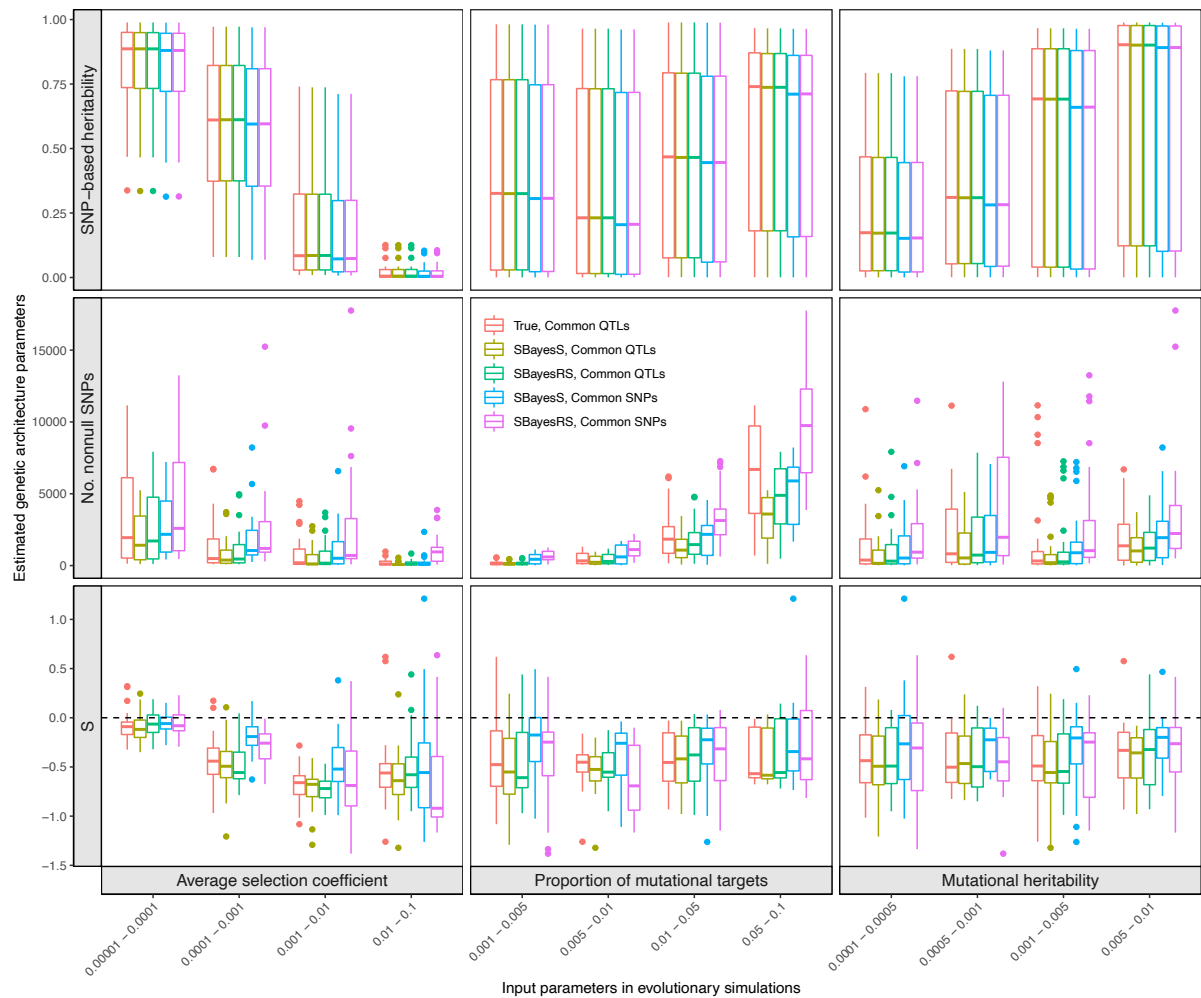
Supplementary Figure 14 Variational patterns of the estimated genetic architecture parameters under different scenarios of evolutionary simulations, when selection coefficients followed a mixture distribution. The Simons et al. pleiotropic model was used to generate genetic effects. The x-axis shows the values of three input parameters in evolutionary simulations. The y-axis shows the distribution of the genetic architecture parameter estimated by SBayesS. Colour shows the results under the model of Simons et al with $n_t = 1, 2, 4$ or 10 . It can be seen that different n_t only affected the SNP-based heritability. The band inside the box is the median, the bottom and top of the box are the first and third quartiles, respectively (Q1 and Q3), and the lower and upper whiskers are $Q1 - 1.5 \text{ IQR}$ and $Q3 + 1.5 \text{ IQR}$, respectively, where $\text{IQR} = Q3 - Q1$.



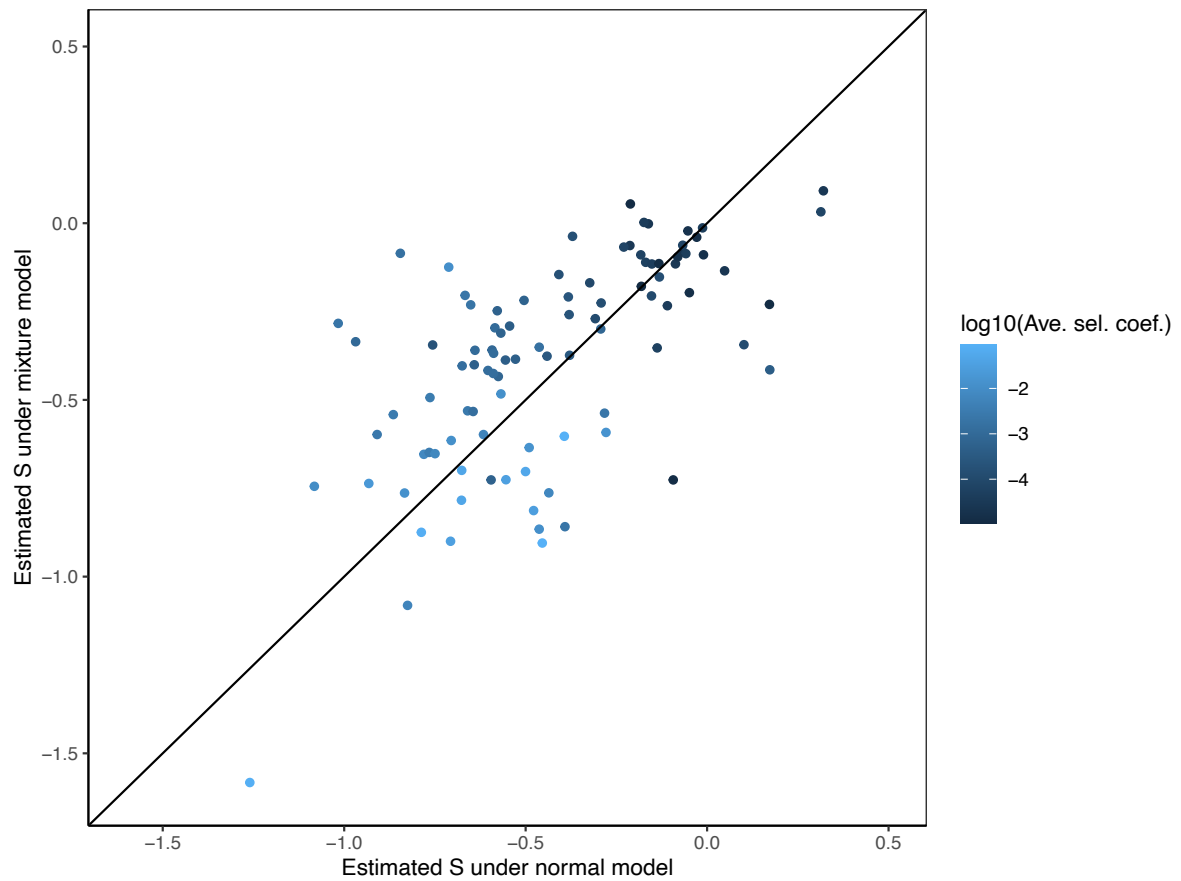
Supplementary Figure 15 Variational patterns of the estimated genetic architecture parameters under different scenarios of evolutionary simulations, when selection coefficients followed a mixture distribution. The Eyre-Walker model was used to generate genetic effects ($\tau = 0.5$ and $\sigma^2 = 0.1$). The x-axis shows the values of three input parameters in evolutionary simulations. The y-axis shows the distribution of the genetic architecture parameter estimates, where the polygenicity parameter is represented by the number of nonnull SNPs for better benchmarking. “True, Common QTLs”: parameters computed directly from the simulated genetic effects of all common causal variants; “SBayesS, Common QTLs” (or “SBayesRS, Common QTLs”): SBayesS (or SBayesRS) estimates using the genotype data of the common causal variants and the phenotypes; “SBayesS, Common SNPs” (or “SBayesRS, Common SNPs”): SBayesS (or SBayesRS) estimates using the genotype data of 36k common SNPs and the simulated genetic values. The band inside the box is the median, the bottom and top of the box are the first and third quartiles, respectively (Q1 and Q3), and the lower and upper whiskers are $Q1 - 1.5 \text{ IQR}$ and $Q3 + 1.5 \text{ IQR}$, respectively, where $\text{IQR} = Q3 - Q1$.



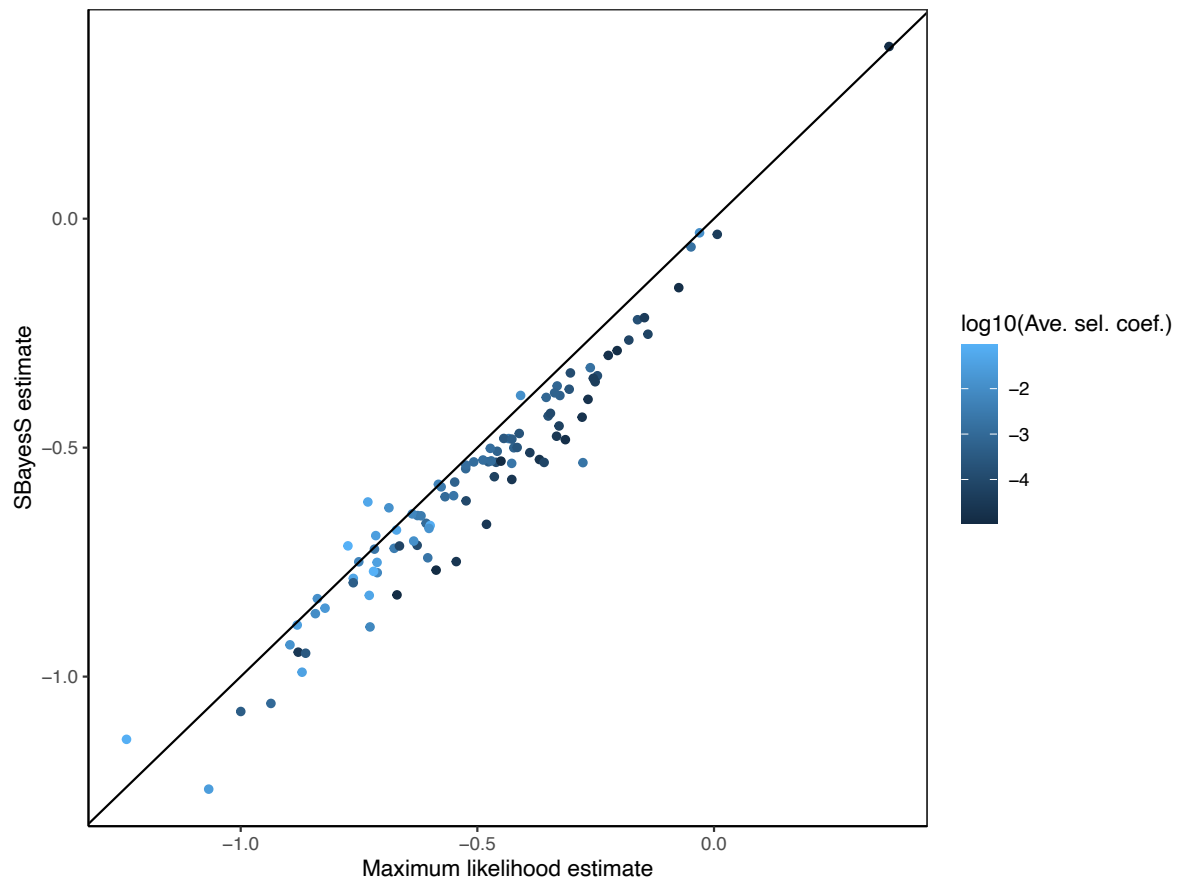
Supplementary Figure 16 Variational patterns of the estimated genetic architecture parameters under different scenarios of evolutionary simulations, when selection coefficients followed a mixture distribution. The Eyre-Walker model was used to generate genetic effects. The x-axis shows the values of three input parameters in evolutionary simulations. The y-axis shows the distribution of the genetic architecture parameters estimated by SBayesS. Colour shows the results under the Eyre-Walker model with $\tau = 0.2, 0.5, 0.8$ or 1 . It can be seen that the genetic architecture parameter estimates were subject to τ . Each box plot shows the results of 25 independent simulation replicates. The band inside the box is the median, the bottom and top of the box are the first and third quartiles, respectively (Q1 and Q3), and the lower and upper whiskers are $Q1 - 1.5 \text{ IQR}$ and $Q3 + 1.5 \text{ IQR}$, respectively, where $\text{IQR} = Q3 - Q1$.



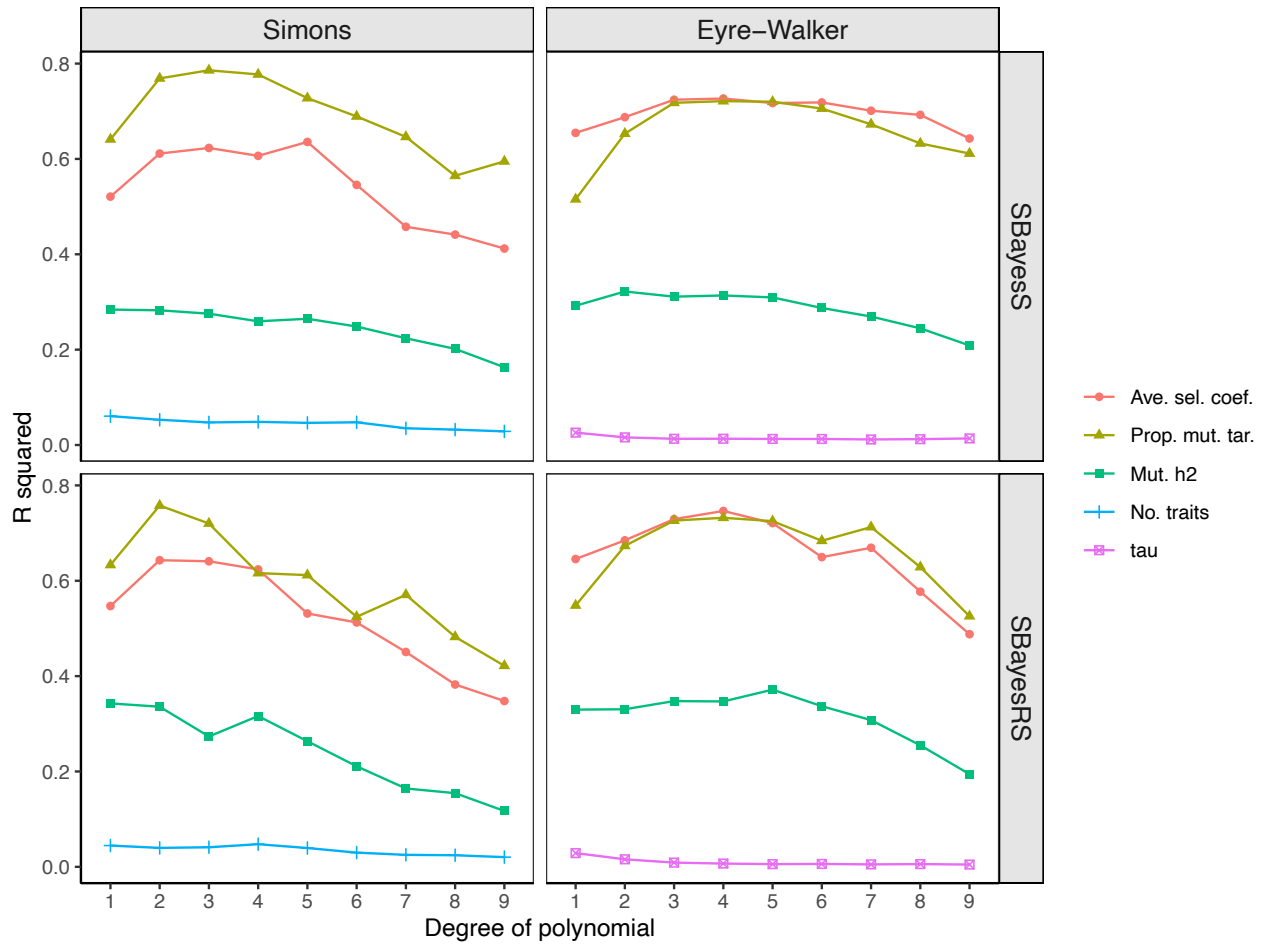
Supplementary Figure 17 Variational patterns of the estimated genetic architecture parameters under different evolutionary simulation scenarios, when selection coefficients followed a normal distribution. The Simons et al. pleiotropic model with $n_t = 1$ was used to generate genetic effects (Methods). The x-axis shows the values of three input parameters in the evolutionary simulations. The y-axis shows the distribution of the genetic architecture parameter estimates, where the polygenicity parameter is represented by the number of nonnull SNPs for better benchmarking. “True, Common QTLs”: parameters computed directly from the simulated genetic effects of all common causal variants; “SBayesS, Common QTLs” (or “SBayesRS, Common QTLs”): SBayesS (or SBayesRS) estimates using the genotype data of the common causal variants and the phenotypes; “SBayesS, Common SNPs” (or “SBayesRS, Common SNPs”): SBayesS (or SBayesRS) estimates using the genotype data of 36k common SNPs and the simulated genetic values. Each box plot shows the results of 25 independent simulation replicates. The band inside the box is the median, the bottom and top of the box are the first and third quartiles, respectively (Q1 and Q3), and the lower and upper whiskers are $Q1 - 1.5 \text{ IQR}$ and $Q3 + 1.5 \text{ IQR}$, respectively, where $\text{IQR} = Q3 - Q1$.



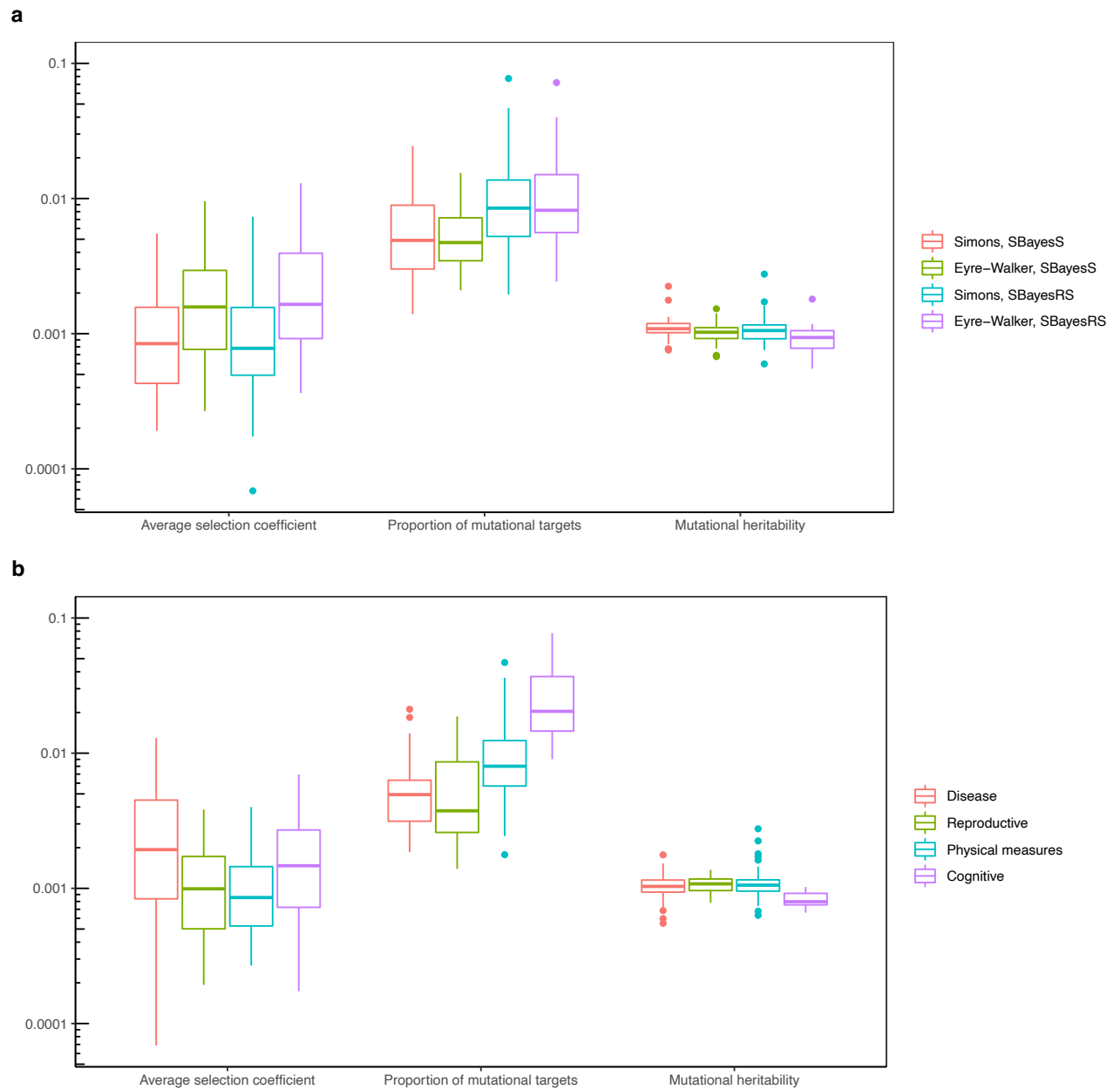
Supplementary Figure 18 Comparison of OLS estimates of S based on the true causal effects under the mixture distribution and the normal distribution for simulating the selection coefficients.



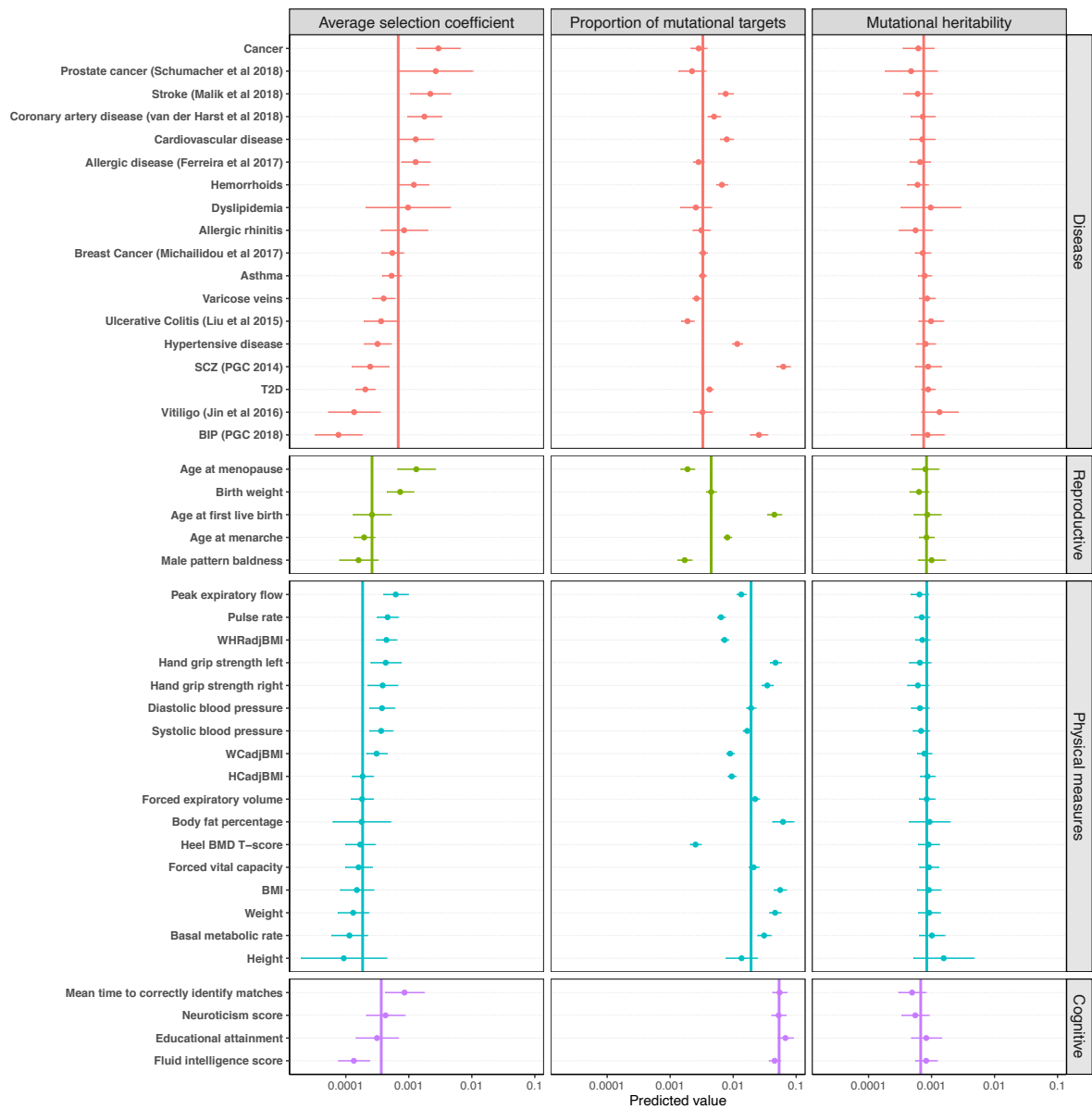
Supplementary Figure 19 SBayesS estimates of S using causal variant genotypes are consistent with the MLE estimates based on the causal variant effects.



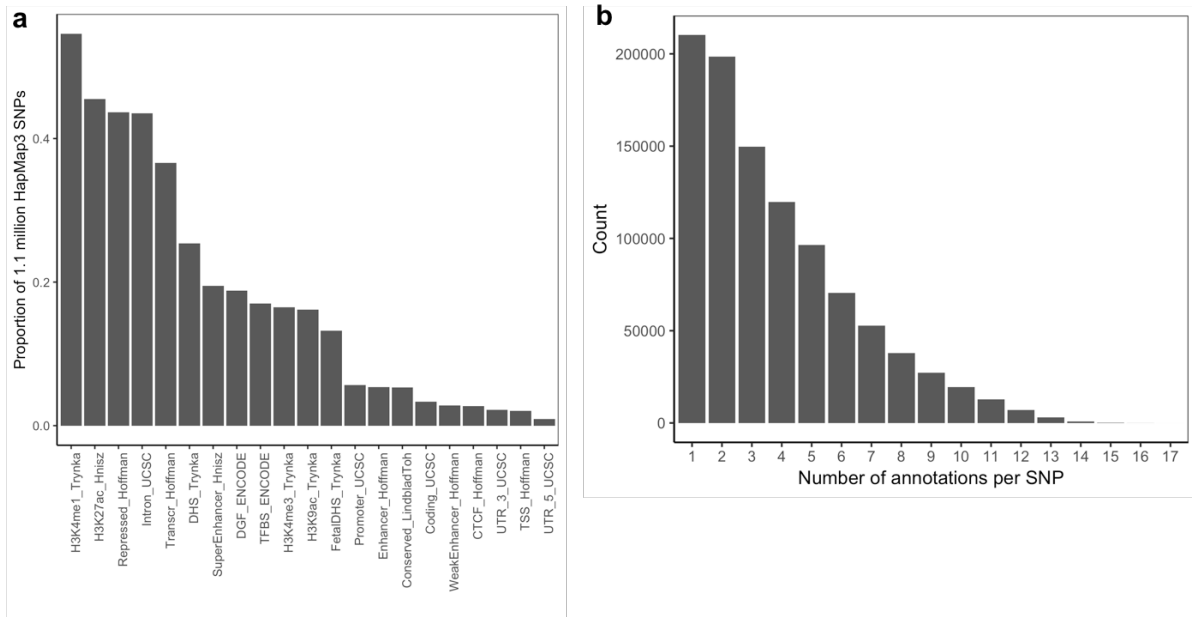
Supplementary Figure 20 Prediction R^2 in cross-validation for predicting the evolutionary parameters using a polynomial regression with the simulated data under the Simons et al or Eyre-Walker model. The predictors in the polynomial regression were the estimated genetic architecture parameters, which were \hat{h}_{SNP}^2 , $\hat{\pi}$ and \hat{S} in SBayesS or \hat{h}_{SNP}^2 , $\hat{\pi}_1$, $\hat{\pi}_2$, $\hat{\pi}_3$, $\hat{\pi}_4$ and \hat{S} in SBayesRS. The response variables were \bar{s} (average selection coefficient), π_m (proportion of mutational targets) or h_m^2 (mutational heritability) at log10 scale, or exclusive parameters specific to each model, namely, the number of traits (n_t) in the model of Simons et al. and τ and σ^2 in the Eyre-Walker's model. The prediction R^2 was computed from cross-validation with 80% of simulation data used as training and the rest as validation. It can be seen that there was reasonably high power to predict \bar{s} , π_m and h_m^2 but no power to predict n_t and τ .



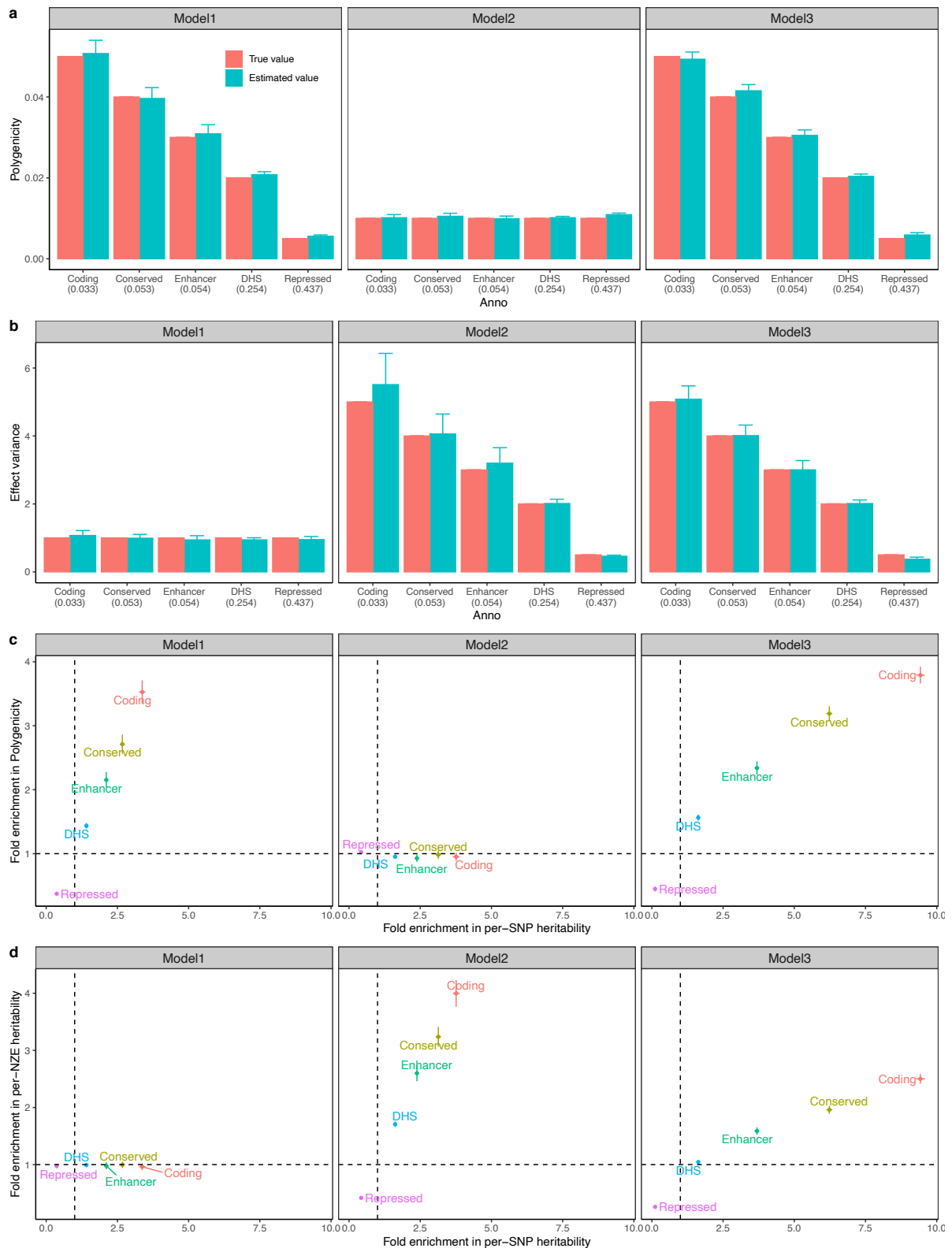
Supplementary Figure 21 Prediction of the evolutionary parameters for 44 complex traits and diseases based on a negative selection model where selection coefficients followed a normal distribution. a) Distribution of the predicted evolutionary parameters under different scenarios: methods used for estimating the genetic architecture parameters (SBayesS and SBayesRS) and pleiotropic effect models used for simulations (the Simons et al. and Eyre-Walker models). b) Distribution of predicted evolutionary parameters for four trait categories. Each box plot shows the results for a number of traits in a category, with each trait having four results from analyses using different estimation methods and simulation models. The band inside the box is the median, the bottom and top of the box are the first and third quartiles, respectively (Q1 and Q3), and the lower and upper whiskers are $Q1 - 1.5 \text{ IQR}$ and $Q3 + 1.5 \text{ IQR}$, respectively, where $\text{IQR} = Q3 - Q1$.



Supplementary Figure 22 Predicted values of the evolutionary parameters based on the Simons et al. model for simulation and SBayesS estimates of genetic architecture parameters. The horizontal bar shows the standard error of the predicted value. The vertical bar shows the median value in each trait category.

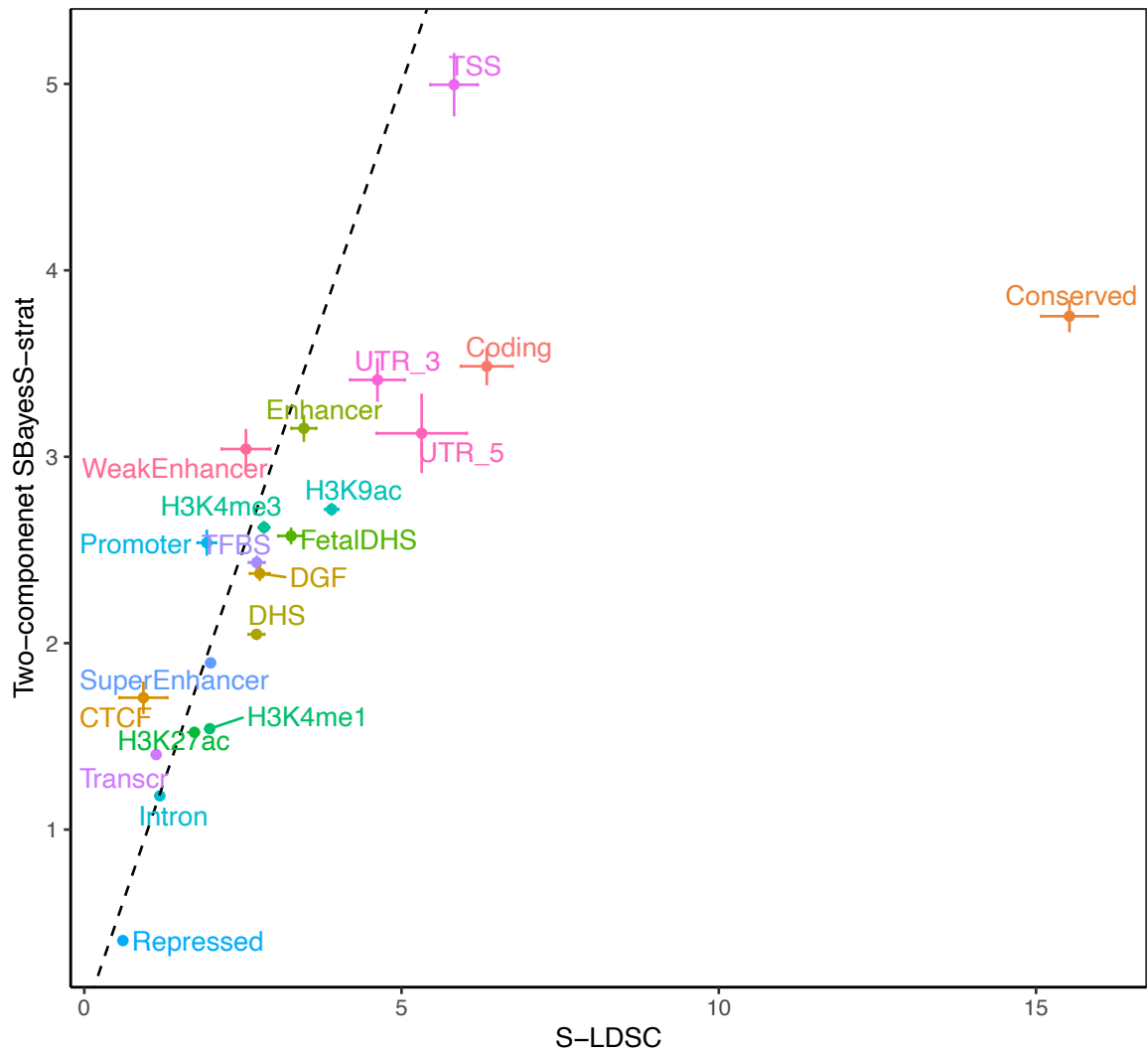


Supplementary Figure 23 Summary of the 21 functional annotation categories from the LDSC baseline model. a) The proportion of 1.1 million HapMap3 common SNPs used in the analysis in each functional category. b) The distribution of the number of annotations for each SNP.

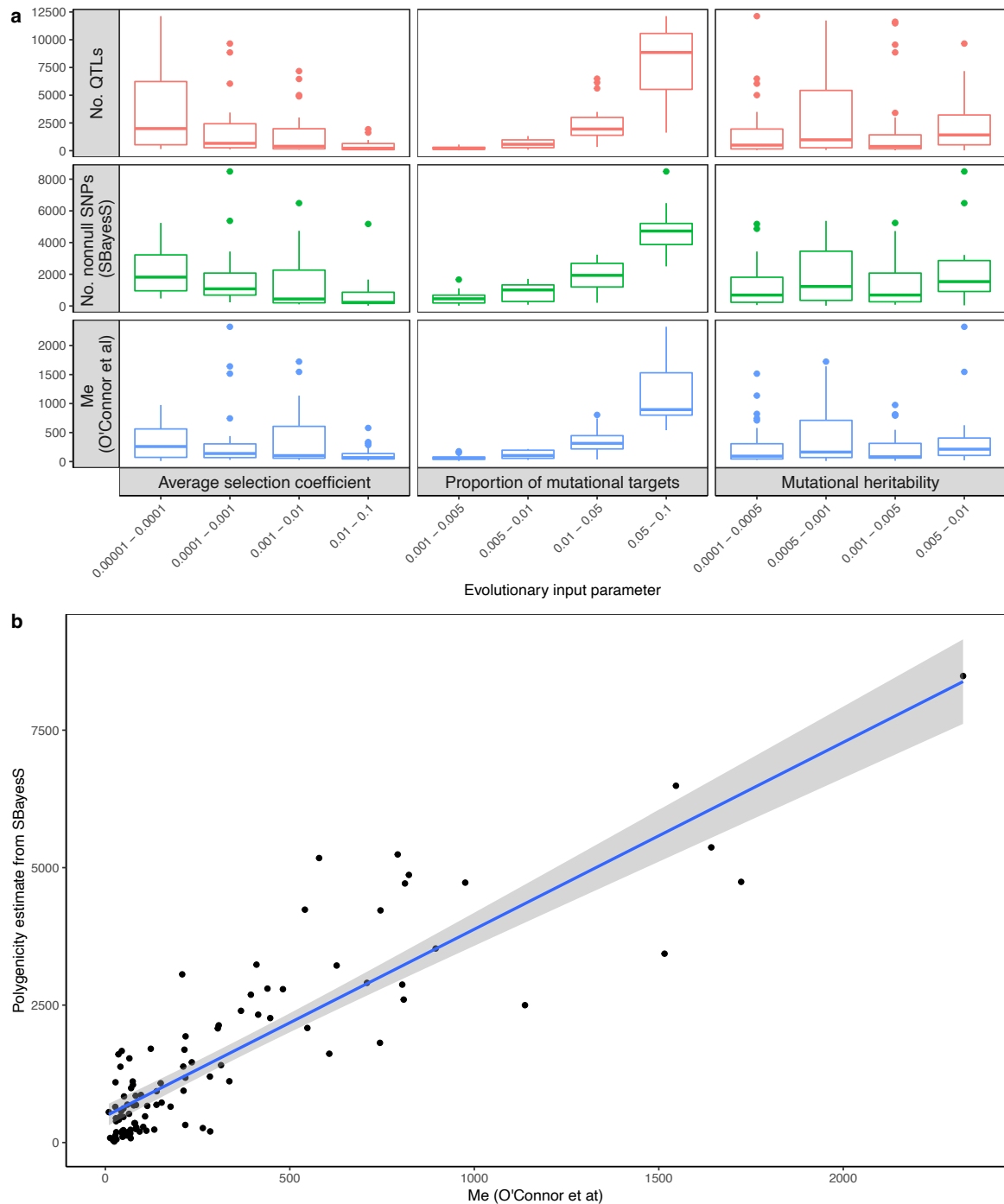


Supplementary Figure 24 Estimation of parameters and their enrichment under different simulated models of genetic architecture across functional annotation categories. We used the UKB data with ~1.1 million HapMap3 common SNPs for simulation and considered three models to simulate the distribution of causal effects. In model 1, the proportion of causal variants was set to 0.5, 2, 3, 4 and 5% for the repressed, DHS, enhancer, conserved and coding regions, respectively, and 1% for the rest of the genome, with the variance of causal effects = 1.

In model 2, the variance of causal effects was set to 0.5, 2, 3, 4 and 5 for the five functional annotation categories and 1 for the rest of the genome, with the proportion of causal variants = 1%. In model 3, we varied both the proportion of causal variants and the variance of causal effects across categories as in model 1 and 2. The trait heritability was set to 0.5. The annotation data were selected from the LDSC baseline model, which had extensive annotation overlaps. When a causal variant had more than one annotation, we randomly assigned one of the overlapping annotations to the causal variant in the simulation. In the analysis, we used the true annotations, where all annotations were mutually disjoint in each replicate, and then run the 2-component SBayesS-strat model. Results are a) Estimated polygenicity parameter under different models; b) Estimated effect variance under in different models; c) Fold enrichment in polygenicity compared with that in per-SNP heritability under different models; and d) Fold enrichment in per-NZE heritability compared with that in per-SNP heritability under different models. In a) and b), the true value is shown in red bar and the number in the a-axis label shows the proportion of SNPs in the corresponding annotation. Data are presented as mean values +/- standard errors of the means across 30 independent simulation replicates.



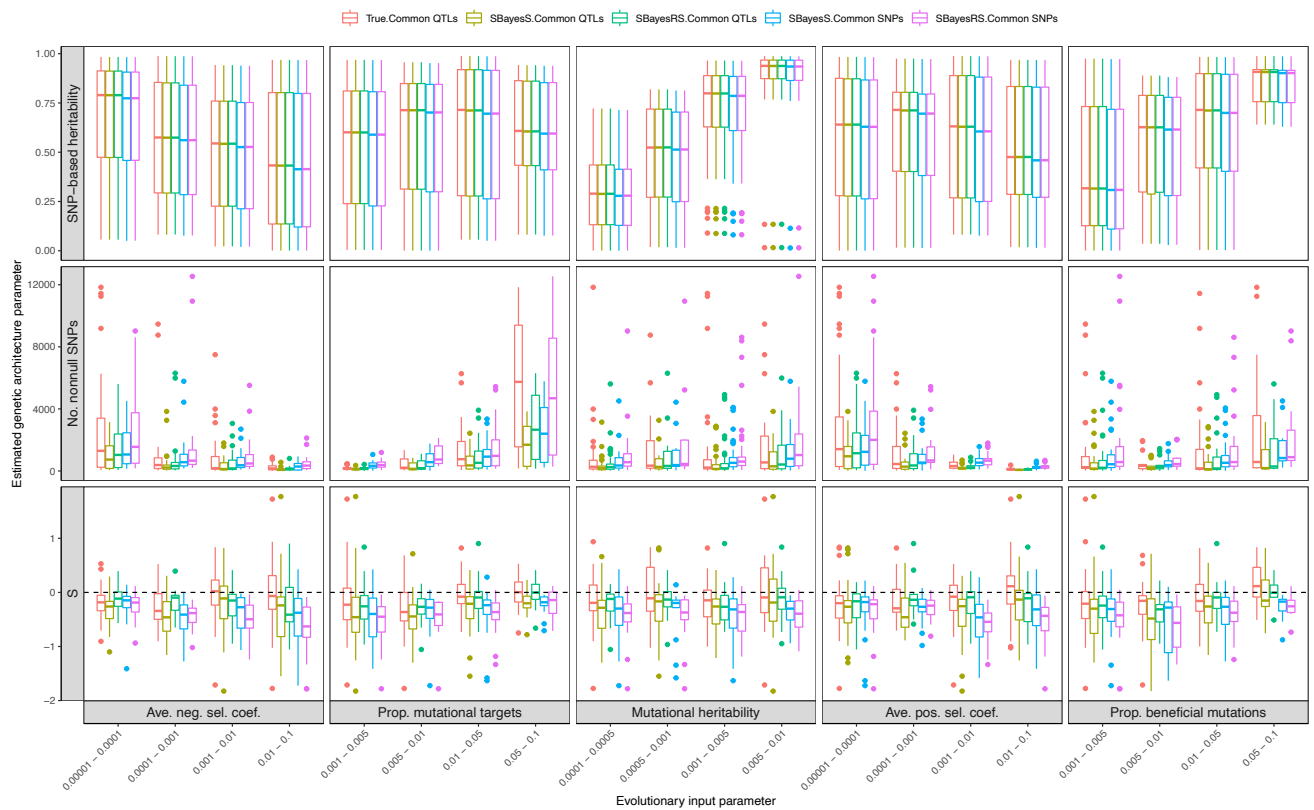
Supplementary Figure 25 Comparison between SBayesS-strat and S-LDSC in per-SNP heritability enrichment using 21 annotation categories from LDSC baseline model. Results are SBayesS-strat analysis that fitted only two components (SNPs in one annotation as the first component and the other SNPs as the second component) versus S-LDSC fitted all annotations. We ran S-LDSC with the same annotations and GWAS summary statistics as used in the analysis above and LD data from the 1000 Genomes Project (the default setting of S-LDSC). Each bar indicates the standard error of the mean. The dashed line shows $y=x$.



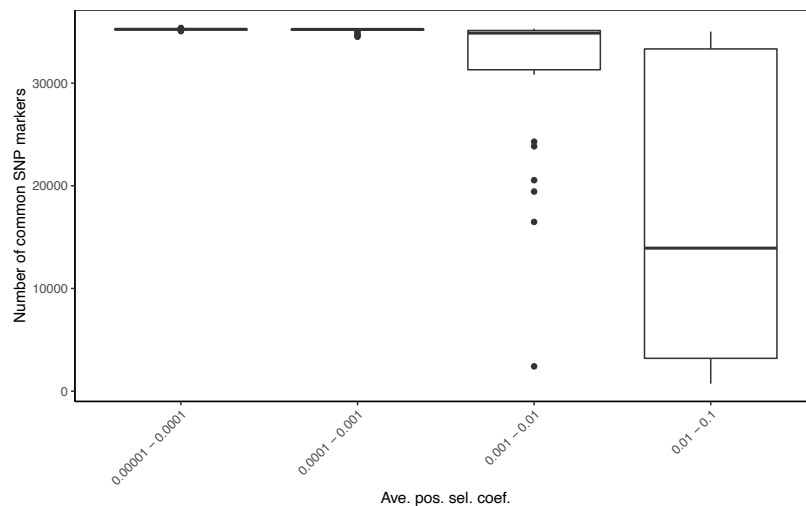
Supplementary Figure 26 Comparison of the polygenicity parameter defined in our study (π) and that in O'Connor et al (M_e) in our forward simulations under negative selection. a) Both π and M_e changed with the total number of causal variants in the simulation. According to the definition in O'Connor et al, $M_e = 3M/\kappa$, $\kappa = E[\beta^4]/E[\beta^2]^2$, where M is the total number of variants (causal variants + 36k SNP markers) and β are the true effect size for the causal variants in per-normalized-genotype units and zero for SNP markers. Each box plot shows the results of 25 independent simulation replicates. The band inside the box is the median, the bottom and top of the box are the first and third quartiles, respectively (Q1 and Q3), and the lower and upper whiskers are $Q1 - 1.5 \text{ IQR}$ and $Q3 + 1.5 \text{ IQR}$, respectively, where $\text{IQR} = Q3 - Q1$.

b) M_e was highly correlated with the polygenicity estimate from SBayesS using 36k common SNP markers (correlation = 0.876, slope of regressing SBayesS estimates on $M_e = 3.4$). The grey bands around the line represent the standard error of the regression line.

a

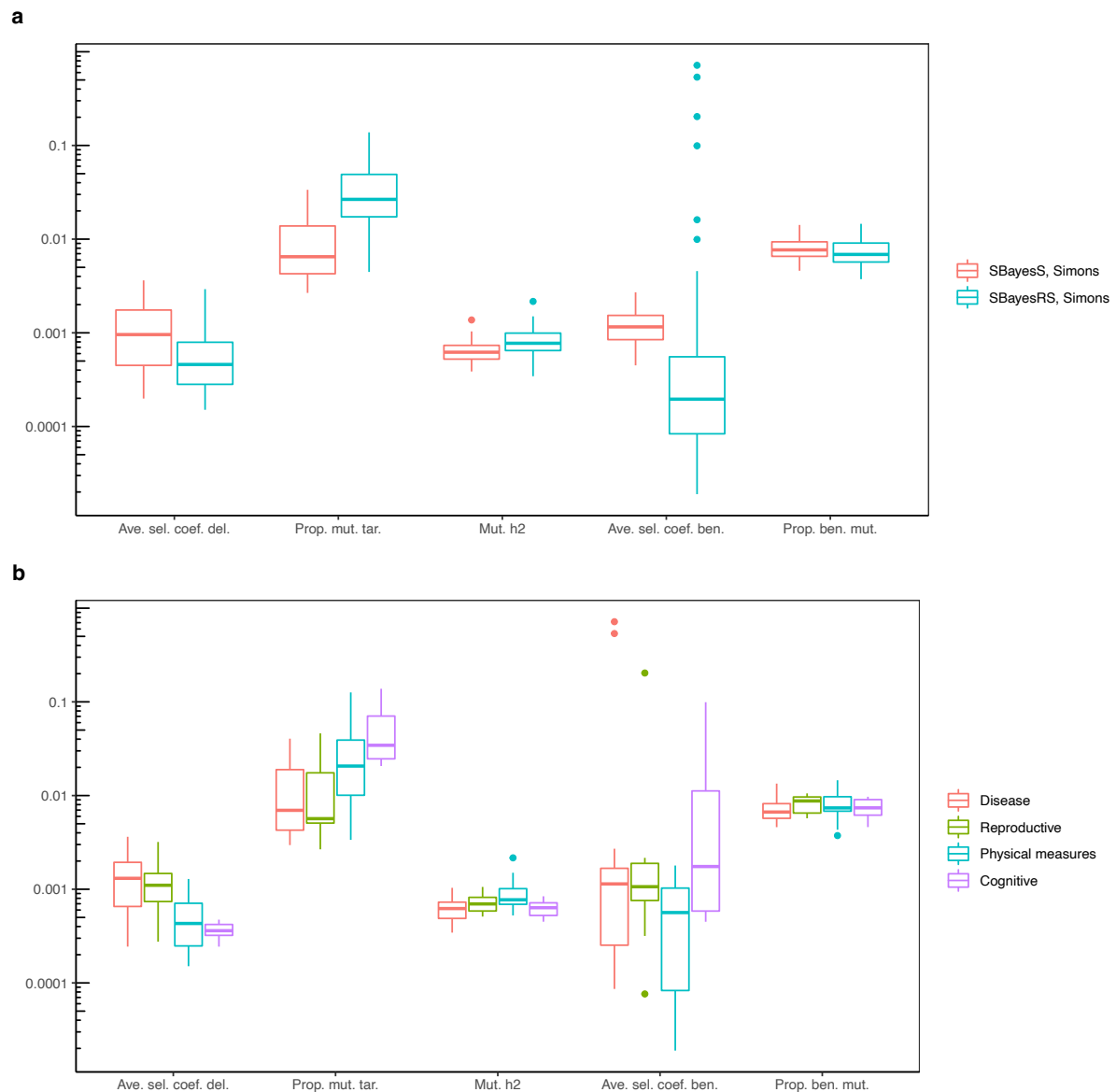


b

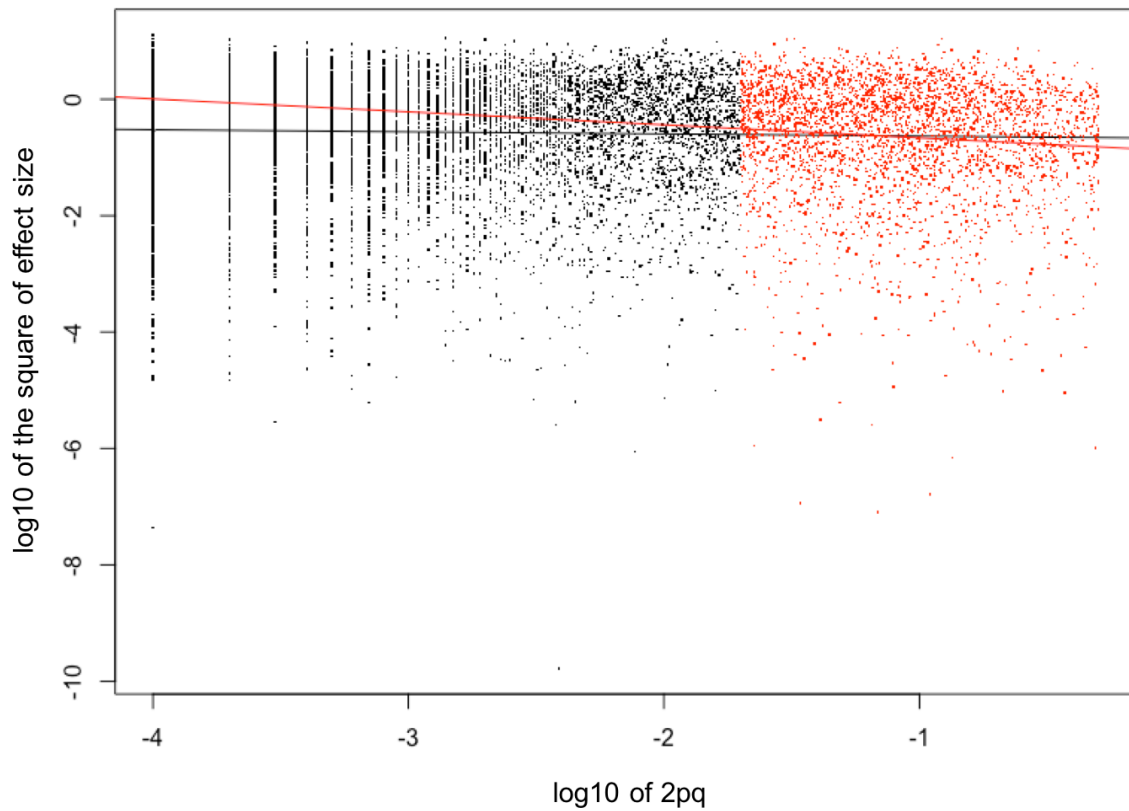


Supplementary Figure 27 Variational patterns of the estimated genetic architecture parameters under both negative and positive selection in evolutionary simulations, when selection coefficients followed a mixture distribution. a) The Simons et al. pleiotropic model with $n_t = 1$ was used to generate genetic effects. The x-axis shows the values of five input parameters in evolutionary simulations. The y-axis shows the distribution of the genetic architecture parameter estimates, where the polygenicity parameter is represented by the number of nonnull SNPs for better benchmarking. “True, Common QTLs”: parameters computed directly from the simulated genetic effects of all common causal variants; “SBayesS, Common

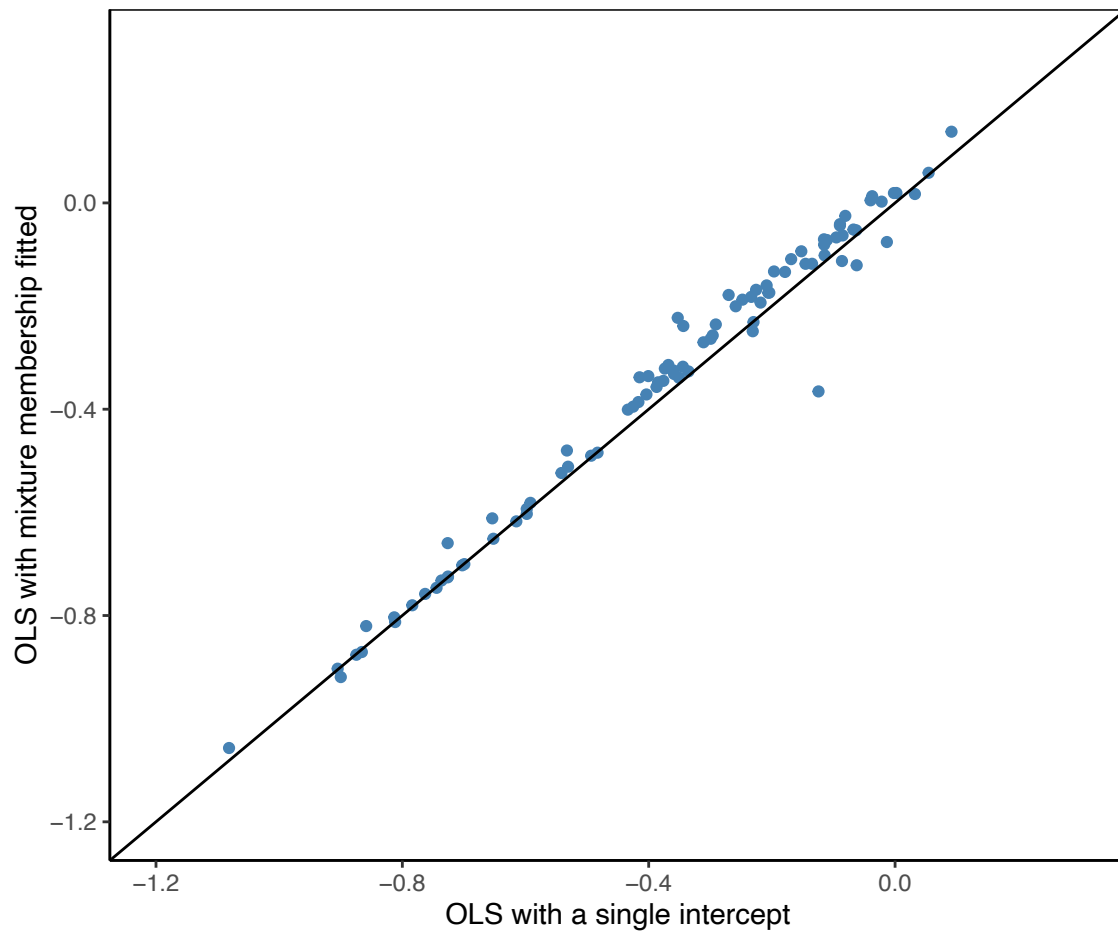
QTLs" (or "SBayesRS, Common QTLs"): SBayesS (or SBayesRS) estimates using the genotype data of the common causal variants and the phenotypes; "SBayesS, Common SNPs" (or "SBayesRS, Common SNPs"): SBayesS (or SBayesRS) estimates using the genotype data of 36k common SNPs and the simulated genetic values. b) The number of common SNP markers decreased with the increased strength of positive selection. Each box plot shows the results of 25 independent simulation replicates. The band inside the box is the median, the bottom and top of the box are the first and third quartiles, respectively (Q1 and Q3), and the lower and upper whiskers are $Q1 - 1.5 \text{ IQR}$ and $Q3 + 1.5 \text{ IQR}$, respectively, where $\text{IQR} = Q3 - Q1$. An apparent discrepancy between the true S value and the estimate from SBayesS/RS using SNP markers was observed when the strength of positive selection increased (row 3, column 4 in panel a). This is likely because the strong positive selection introduced a selective sweep on SNPs in LD with the beneficial mutations of relatively large effects (known as hitchhiking), resulting in a reduced number of common SNPs and a reduced capacity of SNPs to track all the beneficial mutations. In contrast, we did not observe a reduction in the number of common SNPs in the presence of strong negative selection, suggesting background selection has a smaller impact on reducing the SNP diversity than hitchhiking so that the deleterious mutations are better tracked by the SNPs.



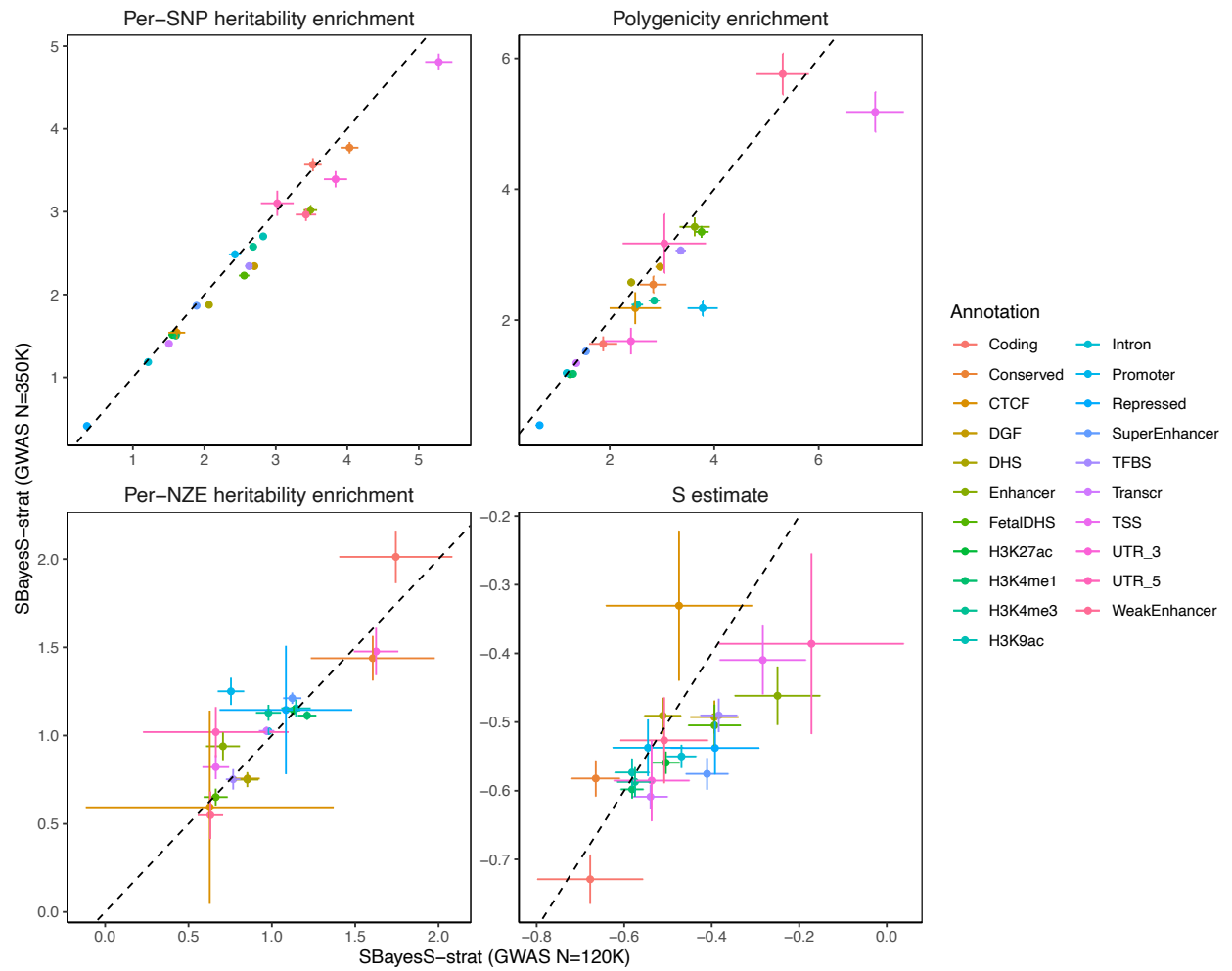
Supplementary Figure 28 Prediction of the evolutionary parameters for 44 complex traits and diseases based on a mixture of negative and positive selection model. a) Distribution of the predicted evolutionary parameters under different methods for estimating genetic architecture parameters (SBayesS and SBayesRS), shown by colours. Each box plot shows the results for 44 complex traits. b) Distribution of predicted evolutionary parameters for five trait categories, shown by colours. Each box plot shows the results for a number of traits in a category, with each trait having two results from analyses using different estimation methods. The band inside the box is the median, the bottom and top of the box are the first and third quartiles, respectively (Q1 and Q3), and the lower and upper whiskers are $Q1 - 1.5 \text{ IQR}$ and $Q3 + 1.5 \text{ IQR}$, respectively, where $\text{IQR} = Q3 - Q1$.



Supplementary Figure 29 Joint distribution of the squared effect size and heterozygosity ($2pq$) for causal variants in the forward simulation in the presence of negative selection. Red colour shows the causal variants with $MAF > 0.01$. The line is the regression line for all causal variants (black) or common causal variants (red), which is an estimate of the S parameter.



Supplementary Figure 30 Ordinary least squares (OLS) estimate of S from a regression model fitting mixture membership-specific intercepts vs. that from a model fitting a single intercept, when the true model is a mixture distribution.



Supplementary Figure 31 Comparison in genetic architecture estimation in 21 functional categories between SBayesS-strat with GWAS sample sizes of 350k and 120k using 1.1M common HapMap3 SNPs. Data are presented as mean values +/- standard errors of the means.

Supplementary Tables

Supplementary Table 1: Estimation of genetic architecture parameters for 35 complex traits (including diseases) in the full UKB data.

The point estimate is the posterior mean and SE is the posterior standard error estimated from the MCMC sample.

The column of "Sample size" for a disease shows the number of cases. The total number of cases and controls is 452,272 for all of the diseases.

Category	Trait	Acronym	Sample size	SNP-based heritability		Polygenicity		S	
				Estimate	SE	Estimate	SE	Estimate	SE
Physical measures	Basal metabolic rate	BMR	339237	0.3143	0.0026	0.0311	0.0009	-0.6209	0.0256
Physical measures	BMI	BMI	344306	0.2619	0.0024	0.0449	0.0016	-0.5461	0.0303
Physical measures	Body fat percentage	BFP	339098	0.2514	0.0024	0.0503	0.0020	-0.5533	0.0308
Physical measures	Hand grip strength left	HGSL	344014	0.1315	0.0021	0.0335	0.0019	-0.5613	0.0448
Physical measures	Hand grip strength right	HGSR	344033	0.1316	0.0021	0.0287	0.0015	-0.4982	0.0451
Physical measures	HCadjBMI	HC	344252	0.2307	0.0025	0.0158	0.0005	-0.5913	0.0294
Physical measures	Heel BMD T-score	HBMD	196375	0.2553	0.0036	0.0058	0.0003	-0.4959	0.0441
Physical measures	Height	Height	344664	0.5452	0.0028	0.0201	0.0004	-0.6529	0.0220
Physical measures	WCadjBMI	WC	344248	0.1734	0.0023	0.0143	0.0006	-0.5899	0.0347
Physical measures	Weight	Weight	344434	0.2763	0.0025	0.0391	0.0013	-0.5620	0.0288
Physical measures	WHRadjBMI	WHR	344228	0.1426	0.0022	0.0120	0.0005	-0.5779	0.0385
Physical measures	Diastolic blood pressure	DBP	322683	0.1410	0.0023	0.0212	0.0010	-0.5417	0.0415
Physical measures	Forced expiratory volume in 1 second	FEV	315184	0.2249	0.0026	0.0249	0.0009	-0.6020	0.0318
Physical measures	Forced vital capacity	FVC	315012	0.2504	0.0026	0.0246	0.0009	-0.6309	0.0287
Physical measures	Peak expiratory flow	PEF	315184	0.1107	0.0022	0.0166	0.0010	-0.5846	0.0475
Physical measures	Pulse rate	PR	325075	0.1384	0.0023	0.0109	0.0005	-0.5578	0.0409
Physical measures	Systolic blood pressure	SBP	322679	0.1471	0.0023	0.0198	0.0009	-0.5603	0.0395
Cognitive	Educational attainment	EA	326945	0.1820	0.0024	0.0478	0.0025	-0.5960	0.0988
Cognitive	Fluid intelligence score	FIS	113198	0.2540	0.0054	0.0384	0.0039	-0.4990	0.0364
Cognitive	Mean time to correctly identify matches	MTCIM	342712	0.0723	0.0019	0.0332	0.0037	-0.4514	0.0738
Cognitive	Neuroticism score	NS	279979	0.1181	0.0023	0.0357	0.0030	-0.4236	0.0604
Reproductive	Age at first live birth	AFLB	125831	0.1966	0.0049	0.0353	0.0043	-0.6681	0.0684
Reproductive	Age at menopause	Mnps	106965	0.0961	0.0049	0.0014	0.0002	-0.6907	0.0818
Reproductive	Age menarche	Mnrch	181335	0.2208	0.0039	0.0143	0.0008	-0.5647	0.0439
Reproductive	Birth weight	BW	196388	0.1046	0.0032	0.0075	0.0006	-0.5175	0.0650
Reproductive	Male pattern baldness	MPB	158696	0.2895	0.0036	0.0035	0.0002	-0.5222	0.0499
Disease	Allergic rhinitis	AR	28041	0.0894	0.0048	0.0049	0.0006	-0.4070	0.0998
Disease	Asthma	Asthma	58479	0.1421	0.0037	0.0060	0.0004	-0.6000	0.0569
Disease	Cancer	Cancer	65534	0.0392	0.0031	0.0027	0.0005	-0.6369	0.0959
Disease	Cardiovascular disease	CD	73856	0.0832	0.0030	0.0110	0.0010	-0.7246	0.0702
Disease	Dyslipidemia	Dyslpl	78921	0.1297	0.0036	0.0038	0.0003	-0.8236	0.0477
Disease	Hemorrhoids	Hmrr	26955	0.0738	0.0047	0.0095	0.0017	-0.5752	0.1227
Disease	Hypertensive disease	HD	87650	0.1757	0.0033	0.0167	0.0009	-0.6460	0.0411
Disease	Type 2 diabetes	T2D	27091	0.2291	0.0060	0.0091	0.0006	-0.5642	0.0507
Disease	Varicose veins	VV	13252	0.1722	0.0082	0.0050	0.0006	-0.6059	0.0794

Supplementary Table 2: Estimation of genetic architecture parameters for 9 common diseases from published GWAS data.
 The point estimate is the posterior mean and SE is the posterior standard error estimated from the MCMC sample.

Disease	Acronym	GWAS summary data set	Cases	Controls	Sample prevalence	Population prevalence	SNP-based heritability		Polygenicity		S	
							Estimate	SE	Estimate	SE	Estimate	SE
Allergic Disease	AD	Ferreira et al (2017 NG)	180129	180709	0.499	0.2	0.0803	0.0026	0.0037	0.0003	-0.5786	0.0772
Bipolar Disorder	BIP	PGC (2018 Cell)	20129	54065	0.271	0.03	0.3439	0.0088	0.0290	0.0050	-0.4198	0.1079
Breast Cancer	BC	Michailidou et al (2017 Nature)	122977	105974	0.537	0.07	0.1338	0.0023	0.0060	0.0003	-0.5553	0.0487
Coronary Artery Disease	CAD	van der Harst et al (2018 Circ Res)	122733	424528	0.224	0.07	0.0713	0.0013	0.0071	0.0004	-0.7452	0.0543
Prostate Cancer	PC	Schumacher et al (2018 NG)	79194	61112	0.564	0.0012	0.0284	0.0010	0.0012	0.0001	-0.3946	0.1053
Schizophrenia	SCZ	PGC (2014 Nature)	36989	113075	0.246	0.01	0.2100	0.0026	0.0462	0.0030	-0.6130	0.0445
Stroke	Stroke	Malik et al (2018 NG)	40585	406111	0.091	0.05	0.0470	0.0026	0.0097	0.0016	-0.6664	0.1355
Ulcerative Colitis	UC	Liu et al (2015 NG)	6968	20464	0.254	0.005	0.1983	0.0104	0.0031	0.0004	-0.6619	0.0904
Vitiligo	Vtlg	Jin et al (2016 NG)	4680	39586	0.106	0.002	0.3818	0.0113	0.0079	0.0006	-0.6960	0.0678

Supplementary Table 3: Classification of 35 UKB traits (including diseases) based on the information provided by the UKB.

Category	Trait	UKB Data-Field code	UKB Category description	UKB Category ID	Remark
Physical measures	Basal metabolic rate	23105	Impedance measures - Anthropometry - Physical measures	100009	
Physical measures	BMI	21001	Body size measures - Anthropometry - Physical measures	100010	
Physical measures	Body fat percentage	23099	Impedance measures - Anthropometry - Physical measures	100009	
Physical measures	Hand grip strength left	46	Hand grip strength - Physical measures	100019	
Physical measures	Hand grip strength right	47	Hand grip strength - Physical measures	100019	
Physical measures	HCadjBMI	49	Body size measures - Anthropometry - Physical measures	100010	Adjusted the phenotype for
Physical measures	Heel BMD T-score	78	Bone-densitometry of heel - Physical measures	100018	
Physical measures	Height	50	Body size measures - Anthropometry - Physical measures	100010	
Physical measures	WCadjBMI	48	Body size measures - Anthropometry - Physical measures	100010	Adjusted the phenotype for
Physical measures	Weight	21002	Body size measures - Anthropometry - Physical measures	100010	
Physical measures	WHRadjBMI	-	-	-	Derived from WC and HC and adjusted for BMI
Physical measures	Diastolic blood pressure	4079	Blood pressure - Physical measures	100011	
Physical measures	Forced expiratory volume in 1 second	3063	Spirometry - Physical measures	100020	
Physical measures	Forced vital capacity	3062	Spirometry - Physical measures	100020	
Physical measures	Peak expiratory flow	3064	Spirometry - Physical measures	100020	
Physical measures	Pulse rate	102	Blood pressure - Physical measures	100011	
Physical measures	Systolic blood pressure	4080	Blood pressure - Physical measures	100011	
Cognitive	Educational attainment	6138	Education - Sociodemographics - Touchscreen	100063	
Cognitive	Fluid intelligence score	20016	Fluid intelligence / reasoning - Cognitive function	100027	
Cognitive	Mean time to correctly identify matches	20023	Reaction time - Cognitive function	100032	
Cognitive	Neuroticism score	20127	Mental health - Psychosocial factors - Touchscreen	100060	
Reproductive	Age at first live birth	2754	Female-specific factors - Sex-specific factors - Touchscreen	708	
Reproductive	Age at menopause	3581	Female-specific factors - Sex-specific factors - Touchscreen	708	
Reproductive	Age at menarche	2714	Female-specific factors - Sex-specific factors - Touchscreen	708	
Reproductive	Birth weight	20022	Early life factors - Verbal interview	708	
Reproductive	Male pattern baldness	2395	Male-specific factors - Sex-specific factors - Touchscreen	708	
Disease	Allergic rhinitis	J30 + 1387	Diagnoses - main ICD10 + secondary ICD10 + self-reported	41202+41204+20002	
Disease	Asthma	J45 + 1111	Diagnoses - main ICD10 + secondary ICD10 + self-reported	41202+41204+20002	
Disease	Cancer	C01, C03, C04, C05, C06, C07, C08, C09, C10, C11, C12, C13, C14, C15, C16, C18, C22, C23, C24, C25, C33, C43, C44, C45, C46, C47, C48, C49, C50, C53, C54, C55, C56, C57, C58, C59, C60, C61, C64, C67, C71, C72, C76, C77, C78, C79, C80, C81, C82, C83, C84, C85, C91, C92, C94, C96, D48, K22.7	Diagnoses - main ICD10 + secondary ICD10	41202+41204	Phenotype are acquired from self-reported, ICD 10 main diagnosis and ICD 10 secondary diagnosis
Disease	Cardiovascular disease	I01, I09, I11, I13, I20, I21, I22, I23, I24, I25, I42, I43, I44, I45, I46, I47, I48, I49, I50, I60, I61, I62, I63, I64, I65, I66, I67, I68, I69, I71 + 1074, 1075, 1066, 1485, 1076, 1471, 1483, 1484, 1486, 1487, 1588, 1426, 1479, 1086, 1491, 1083, 1082, 1425, 1592, 1068, 1094, 1079, 1492, 1591, 1583, 1077	Diagnoses - main ICD10 + secondary ICD10 + self-reported	41202+41204+20002	
Disease	Dyslipidemia	E78 + 1473	Diagnoses - main ICD10 + secondary ICD10 + self-reported	41202+41204+20002	
Disease	Hemorrhoids	I84 + 1505	Diagnoses - main ICD10 + secondary ICD10 + self-reported	41202+41204+20002	
Disease	Hypertensive disease	I10-I13 + 1065, 1072, 1073	Diagnoses - main ICD10 + secondary ICD10 + self-reported	41202+41204+20002	
Disease	Type 2 diabetes	E11 + 1223	Diagnoses - main ICD10 + secondary ICD10 + self-reported	41202+41204+20002	
Disease	Varicose veins	I26, I70, I82 + 1067, 1087, 1088	Diagnoses - main ICD10 + secondary ICD10 + self-reported	41202+41204+20002	

Supplementary Table 5: Estimation of genetic architecture parameters for 21 functional genomic categories from the LDSC baseline model meta-analysed over 44 UKB complex traits and diseases.

Results are from the two-component model where the SNPs in one annotation are fitted as one group and the other SNPs are fitted as the other group. The point estimate is the posterior median and SE is the posterior standard error estimated from the MCMC sample. Polygenicity is the proportion of SNPs with nonzero effects among all SNPs in the annotation.

Annotation	Number of SNPs	Fraction of SNPs	SNP-based heritability		Polygenicity		S	
			median	s.e.m.	median	s.e.m.	median	s.e.m.
Coding_UCSC	33338	0.033	0.1158	0.0034	0.0513	0.0031	-0.7392	0.0411
Conserved_LindbladToh	53319	0.053	0.1993	0.0045	0.0686	0.0044	-0.5741	0.0390
CTCF_Hoffman	27260	0.027	0.0464	0.0024	0.0923	0.0075	-0.4981	0.0950
DGF_ENCODE	188706	0.188	0.4463	0.0074	0.0720	0.0024	-0.4983	0.0272
DHS_Trynka	254685	0.254	0.5195	0.0066	0.0669	0.0021	-0.5100	0.0254
Enhancer_Hoffman	54119	0.054	0.1700	0.0039	0.0803	0.0044	-0.4677	0.0391
FetalDHS_Trynka	132853	0.132	0.3408	0.0058	0.0973	0.0034	-0.4594	0.0336
H3K27ac_Hnisz	456432	0.455	0.6923	0.0052	0.0293	0.0005	-0.5646	0.0176
H3K4me1_Trynka	548165	0.546	0.8418	0.0061	0.0333	0.0006	-0.5929	0.0155
H3K4me3_Trynka	165655	0.165	0.4328	0.0057	0.0609	0.0016	-0.5805	0.0222
H3K9ac_Trynka	162379	0.162	0.4398	0.0055	0.0563	0.0017	-0.5862	0.0235
Intron_UCSC	436629	0.435	0.5139	0.0047	0.0237	0.0006	-0.5304	0.0193
Promoter_UCSC	56916	0.057	0.1440	0.0039	0.0477	0.0027	-0.5188	0.0409
Repressed_Hoffman	438171	0.437	0.1767	0.0052	0.0097	0.0011	-0.5735	0.0575
SuperEnhancer_Hnisz	195544	0.195	0.3693	0.0040	0.0298	0.0008	-0.5239	0.0262
TFBS_ENCODE	170829	0.170	0.4142	0.0062	0.0662	0.0021	-0.5399	0.0296
Transcr_Hoffman	367588	0.366	0.5136	0.0054	0.0339	0.0010	-0.6069	0.0202
TSS_Hoffman	20812	0.021	0.1036	0.0035	0.1184	0.0059	-0.3609	0.0657
UTR_3_UCSC	22031	0.022	0.0750	0.0026	0.0636	0.0060	-0.5122	0.0549
UTR_5_UCSC	9223	0.009	0.0287	0.0020	0.0977	0.0088	-0.5521	0.1102
WeakEnhancer_Hoffman	28247	0.028	0.0856	0.0030	0.1350	0.0069	-0.5917	0.0540

References

- [1] Jian Yang, Teresa Ferreira, Andrew P Morris, Sarah E Medland, Genetic of Consortium, DIAbetes Consortium, Pamela AF Madden, Andrew C Heath, Nicholas G Martin, Grant W Montgomery, Michael N Weedon, Ruth J Loos, Timothy M Frayling, McCarthy, Mark I, Joel N Hirschhorn, Michael E Goddard, Peter M Visscher. Conditional and joint multiple-SNP analysis of GWAS summary statistics identifies additional variants influencing complex traits. *Nat Genet*, 44(4):369, 2012.
- [2] Xiang Zhu, Matthew Stephens. Bayesian large-scale multiple regression with summary statistics from genome-wide association studies. *Ann Appl Statistics*, 11(3):1561–1592, 2017.
- [3] Jian Zeng, Ronald Vlaming, Yang Wu, Matthew R Robinson, Lloyd-Jones, Luke R, Loic Yengo, Chloe X Yap, Angli Xue, Julia Sidorenko, McRae, Allan F, Joseph E Powell, Grant W Montgomery, Andres Metspalu, Tonu Esko, Greg Gibson, Naomi R Wray, Peter M Visscher, Jian Yang. Signatures of negative selection in the genetic architecture of human complex traits. *Nat Genet*, strona 1, 2018.
- [4] Lloyd-Jones, Luke R, Jian Zeng, Julia Sidorenko, Loic Yengo, Gerhard Moser, Kathryn E Kemper, Huanwei Wang, Zhili Zheng, Reedik Magi, Tonu Esko, Andres Metspalu, Naomi R Wray, Michael E Goddard, Jian Yang, Peter M Visscher. Improved polygenic prediction by bayesian multiple regression on summary statistics. *Biorxiv*, strona 522961, 2019.
- [5] Mario PL Calus. Right-hand-side updating for fast computing of genomic breeding values. 46(1):24, 2014.
- [6] Jian Yang, Jian Zeng, Michael E Goddard, Naomi R Wray, Peter M Visscher. Concepts, estimation and interpretation of SNP-based heritability. *Nat Genet*, 49(9):1304–1310, 2017.
- [7] O’Connor, Luke J, Alkes L Price. Distinguishing genetic correlation from causation across 52 diseases and complex traits. strony 1728–1734, 2018.